

---

# Contextual Observability and Grammar Singularity for Compositional Task Families

---

Manoj Saravanan<sup>1</sup> Rohit Kumar Salla<sup>1</sup> Shrikar Reddy Kota<sup>1</sup>

## Abstract

We study Bayesian meta-learning when tasks are sampled from a latent compositional family rather than treated independently. A task family is modeled as a stochastic grammar over reusable linear modules, and grouped datasets are observed across many tasks. The target is the latent program law together with the shared module library, modulo explicit language symmetries. We prove four results. First, an exact sufficient-statistic reduction turns grouped compositional meta-learning into a finite mixture of matrix-normal laws. Second, local quotient identifiability holds under positive linearized observability and component separation, while in the single-occurrence regime zero contextual observability yields exact non-identifiability. Third, the posterior for the grouped predictive law contracts at the explicit non-i.i.d. rate

$$\delta_{m,n}^2 = \frac{K_{\text{eff}} \log m + rd^2 \log(mn)}{m},$$

and structural contraction occurs at the inverse exponent  $\kappa$ :  $\kappa = 1$  in regular families but  $\kappa = 2$  in an explicit duplicated-module singular family. Fourth, matching minimax lower bounds identify the same hardness parameters. These quantities yield a theory-native benchmark with provable phase transitions in observability, singularity, and anchor exposure.

## 1. Introduction

Modern evaluation is increasingly organized around *families* of related tasks rather than isolated test points. In compositional regimes, however, benchmark scores remain

---

<sup>1</sup>Virginia Polytechnic Institute and State University, Blacksburg, VA, USA. Correspondence to: Manoj Saravanan <manoj663@vt.edu>.

Accepted to the 1st Workshop on Combining Theory and Benchmarks, CTB@ICML 2026, Seoul, South Korea. Copyright 2026 by the author(s).

fundamentally retrospective: they say how a model performed on a fixed suite, but not when it should transfer to *new compositions* of known primitives, when posterior uncertainty about reusable structure should persist, or which benchmark cells are statistically most informative. In such settings the natural statistical object is not a single task but a *task-generating family*. We study that object directly. Each task first draws a latent program  $Z_i$  from a law  $\pi^*$  on an admissible language  $\mathcal{L}$ , then composes a shared module library  $\Theta^*$  into an operator  $A_{Z_i}$ , and finally generates within-task data. The target is therefore

$$\eta^* = (\pi^*, \Theta^*),$$

defined only up to explicit language symmetries. This grouped-data perspective puts transfer, compositional generalization, and uncertainty over latent task identity into one model: the learner must infer both the *distribution over task programs* and the *reusable module library* from many related datasets.

Recent theory isolates important pieces of this picture, but not the whole object. Special autoregressive task families can permit dramatic extrapolation from small support to exponentially larger compositional families (Abedsoltan et al., 2025). Modular solution discovery and connected-support conditions clarify when reusable modules can be identified compositionally (Schug et al., 2024). Coverage and path-ambiguity analyses show that compositional generalization can fail for structural rather than purely optimization-related reasons (Chang et al., 2025). On the statistical side, posterior contraction for non-i.i.d. experiments (Ghosal & van der Vaart, 2007; Ghosal et al., 2000), singular learning theory (Watanabe, 2009), weakly identifiable mixture geometry (Ho & Nguyen, 2016; 2019), and classical finite-mixture identifiability (Teicher, 1963; Yakowitz & Spragins, 1968) provide the right asymptotic tools. To our knowledge, these threads have not yet been assembled into a finite-sample Bayesian theory of *latent task-program learning* that simultaneously treats identifiability, predictive versus structural contraction, singularity, and theorem-native benchmark design.

The statistical hardness of this problem is organized by three quantities:

$$K_{\text{eff}}, \quad \Delta_{\text{ctx}}, \quad \kappa,$$

namely effective grammar complexity, contextual observability, and grammar singularity order. Here  $K_{\text{eff}}$  controls the difficulty of learning the program law,  $\Delta_{\text{ctx}}$  measures whether primitives are distinguishable in the contexts actually generated by the family, and  $\kappa$  is the inverse exponent that converts predictive proximity into structural proximity. Two auxiliary moduli,  $\lambda_{\text{lin}}$  and  $\gamma_{\text{comp}}$ , determine when positive observability upgrades to local quotient identifiability. The main message is that *prediction and structure can contract at different rates*: regular families have  $\kappa = 1$ , while singular families can have  $\kappa > 1$ .

Our first result is an exact sufficient-statistic reduction: conditional on the within-task designs, least-squares task statistics are sufficient, and grouped meta-learning becomes a finite mixture of matrix-normal laws indexed by latent programs. Our second result is an observability theory: zero contextual observability produces exact observational collapse in single-occurrence languages, whereas  $\lambda_{\text{lin}} > 0$  and  $\gamma_{\text{comp}} > 0$  imply local quotient identifiability. Our third result is a predictive-versus-structural contraction theory: the grouped predictive law contracts at rate  $\delta_{m,n}^2 = (K_{\text{eff}} \log m + rd^2 \log(mn))/m$ , and structural contraction occurs at the inverse exponent  $\kappa$ . In an explicit duplicated-module family the first-order geometry cancels, so structure only contracts at  $O_{\Pi}(\delta_{m,n}^{1/2}/\sqrt{n})$ . Our final result is a matching minimax theory, including grammar-law, anchored-module, zero-observability, and singular lower bounds. These same parameters define a theorem-native benchmark whose cells are labeled by  $(K_{\text{eff}}, \Delta_{\text{ctx}}, \lambda_{\text{lin}}, \gamma_{\text{comp}}, \kappa, s_{\text{anc}})$ .

Section 2 introduces stochastic linear task grammars and the quotient parameterization. Section 3 states the main identifiability, contraction, and minimax theorems. Section 4 summarizes the proof architecture. Section 5 presents the theorem-native benchmark and empirical validation. All proofs are deferred to the appendix.

## 2. Stochastic Linear Task Grammars

We work with a finite-language stochastic grammar over reusable linear modules. This is a theorem-friendly specialization of recent compositional task-family models in which tasks are generated by latent programs built from reusable components (Schug et al., 2024; Abedsoltan et al., 2025). The restriction to bounded-depth linear modules is deliberate: it is the smallest setting in which one can simultaneously prove quotient identifiability, predictive posterior contraction, singular structural rates, and matching lower bounds.

**Latent programs and grouped task data.** Fix integers  $r \geq 2$ ,  $d \geq 1$ , and  $L \geq 1$ , let  $\mathcal{Z}_{L,r} := \bigcup_{t=1}^L [r]^t$ , and let  $\mathcal{L} \subseteq \mathcal{Z}_{L,r}$  be a finite nonempty admissible language. A

program  $z = (z_1, \dots, z_t) \in \mathcal{L}$  composes a module library  $\Theta = (A_1, \dots, A_r) \in \mathbb{A}_B^r$ , where  $\mathbb{A}_B := \{A \in \mathbb{R}^{d \times d} : \|A\|_{\text{op}} \leq B\}$ , via

$$A_z := A_{z_t} \cdots A_{z_1}.$$

A stochastic linear task grammar is

$$\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L}) := \mathfrak{P}_{\mathcal{L}}^{\circ} \times \mathbb{A}_B^r,$$

$$\mathfrak{P}_{\mathcal{L}}^{\circ} := \{\pi \in (0, 1)^{\mathcal{L}} : \sum_{z \in \mathcal{L}} \pi_z = 1\}.$$

For each task  $i$ , draw  $Z_i \sim \pi$ , then  $x_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$  and

$$y_{ij} = A_{Z_i} x_{ij} + \varepsilon_{ij},$$

$$\varepsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2 I_d), \quad j = 1, \dots, n.$$

With  $X_i = [x_{i1}, \dots, x_{in}]$ ,  $Y_i = [y_{i1}, \dots, y_{in}]$ , and  $D_i = (X_i, Y_i)$ , the one-task grouped law is

$$Q_{\eta}^{(n)} = \sum_{z \in \mathcal{L}} \pi_z P_{A_z}^{(n)}, \quad (1)$$

and the meta-dataset satisfies  $D_1, \dots, D_m \stackrel{\text{i.i.d.}}{\sim} Q_{\eta}^{(n)}$ . Throughout, the truth is  $\eta^* = (\pi^*, \Theta^*) \in \Xi(B, \mathcal{L})$  with  $\pi_z^* \geq \pi_{\min} > 0$ .

**Quotient parameterization.** A statistical task family is unchanged by relabeling primitive symbols whenever the admissible language is preserved. The relevant symmetry group is

$$\text{Aut}(\mathcal{L}) := \{\sigma \in S_r : \sigma \cdot \mathcal{L} = \mathcal{L}\},$$

acting on parameters by  $(\sigma \cdot \pi)_u = \pi_{\sigma^{-1}u}$  and  $(\sigma \cdot \Theta)_a = A_{\sigma^{-1}(a)}$ . Structural error is measured by the quotient metric

$$d_q(\eta, \eta') := \min_{\sigma \in \text{Aut}(\mathcal{L})} \left( \|\pi - \sigma \cdot \pi'\|_1 + \sum_{a=1}^r \|A_a - (\sigma \cdot \Theta')_a\|_F \right). \quad (2)$$

This quotient removes only explicit symbol symmetries. Any additional equality of grouped laws that is *not* induced by an automorphism is treated as genuine statistical non-identifiability rather than as notational redundancy. This distinction becomes essential in the low-observability regime.

**Proposition 2.1** (Exact reduced experiment). *Assume  $n \geq d + 1$ . Define the task-level least-squares statistic*

$$\widehat{W}_i := Y_i X_i^{\top} (X_i X_i^{\top})^{-1}, \quad i = 1, \dots, m,$$

which is well defined almost surely under the Gaussian design. Then, conditional on  $X_i$  and  $Z_i = z$ ,

$$\widehat{W}_i \sim \text{MN}_{d,d}(A_z, \sigma^2 I_d, (X_i X_i^{\top})^{-1}),$$

where  $\text{MN}_{d,d}(M, U, V)$  denotes the matrix-normal law with mean  $M$ , row covariance  $U$ , and column covariance  $V$ . Consequently, conditional on the designs  $X_{1:m}$ , the reduced observations satisfy

$$\begin{aligned} Z_i &\stackrel{\text{i.i.d.}}{\sim} \pi, & \widehat{W}_i &= A_{Z_i} + \Xi_i, \\ \Xi_i \mid X_{1:m} &\stackrel{\text{ind}}{\sim} \text{MN}_{d,d}(0, \sigma^2 I_d, (X_i X_i^\top)^{-1}). \end{aligned} \quad (3)$$

Moreover, under any prior on  $\eta$ ,

$$\Pi(\cdot \mid X_{1:m}, Y_{1:m}) = \Pi(\cdot \mid X_{1:m}, \widehat{W}_{1:m}) \quad \text{a.s.}$$

Theorem 2.1 is the technical pivot of the paper: Bayesian meta-learning over latent task grammars becomes an exact grouped finite-mixture problem over the composed operators  $\{A_z\}_{z \in \mathcal{L}}$ . All posterior theory and all experiments in the paper are carried out in this reduced experiment.

### 3. Main Results

At a high level, the theory separates into two layers. The first is a *predictive* layer: after exact sufficient-statistic reduction, Bayesian meta-learning over latent task grammars becomes a non-i.i.d. grouped mixture problem, and posterior contraction is governed by the reduced predictive geometry. The second is a *structural* layer: converting predictive concentration into localization of the latent grammar requires a local inverse inequality, and the resulting exponent is regular ( $\kappa = 1$ ) or singular ( $\kappa > 1$ ) depending on the local geometry. The novelty is that both layers are derived for a stochastic *task-family grammar*, not for an ordinary finite mixture or a single compositional function.

Throughout,  $\Pi_m(\cdot)$  denotes the posterior given the grouped data  $(X_i, Y_i)_{i=1}^m$ , equivalently given the reduced experiment  $(X_i, \widehat{W}_i)_{i=1}^m$  by Theorem 2.1. We write

$$\delta_{m,n}^2 := \frac{K_{\text{eff}} \log m + rd^2 \log(mn)}{m}, \quad K_{\text{eff}} := |\mathcal{L}| - 1. \quad (4)$$

**Hardness quantities.** Besides  $K_{\text{eff}}$ , four structural quantities govern the problem:

$$\Delta_{\text{ctx}}(\eta) := \min_{a \neq b} \sum_{(u,v) \in \mathcal{C}_{\mathcal{L}}} \nu_\pi(u, v) \|A_v(A_a - A_b)A_u\|_F^2, \quad (5)$$

$$\begin{aligned} \lambda_{\text{lin}}(\eta) &:= \inf_{\|H\|_{\text{lib}}=1} \|D\Phi_\Theta(H)\|_\pi^2, \\ \Phi(\Theta) &:= (A_z(\Theta))_{z \in \mathcal{L}}, \end{aligned} \quad (6)$$

$$\gamma_{\text{comp}}(\eta) := \min_{z \neq z'} \|A_z - A_{z'}\|_F. \quad (7)$$

The quantity  $\Delta_{\text{ctx}}$  is the exact context-averaged information carried by one-hole substitutions,  $\lambda_{\text{lin}}$  is the weighted

smallest singular value of the Jacobian of the program-operator map, and  $\gamma_{\text{comp}}$  is the minimum separation between composed operators. The *grammar singularity order*  $\kappa \geq 1$  is the smallest local exponent such that  $\bar{h}_X^2(\eta, \eta^*) \gtrsim d_q(\eta, \eta^*)^{2\kappa}$ . Regular families have  $\kappa = 1$ ; singular families have  $\kappa > 1$ .

#### 3.1. Identifiability and predictive contraction

**Theorem 3.1** (Local quotient identifiability). *Suppose*

$$\lambda_{\text{lin}}(\eta^*) > 0 \quad \text{and} \quad \gamma_{\text{comp}}(\eta^*) > 0.$$

*Then there exists  $\rho_{\text{id}} > 0$  such that for every  $n \geq d + 1$ ,*

$$\begin{aligned} d_q(\eta, \eta^*) &\leq \rho_{\text{id}} \quad \text{and} \quad Q_\eta^{(n)} = Q_{\eta^*}^{(n)} \\ &\implies d_q(\eta, \eta^*) = 0. \end{aligned}$$

Thus pairwise visibility and full local invertibility are distinct:  $\Delta_{\text{ctx}}$  measures substitution-level information, whereas  $\lambda_{\text{lin}}$  controls all infinitesimal perturbations of the library.

**Proposition 3.2** (Exact non-identifiability at zero contextual observability). *Assume that  $\mathcal{L}$  is single-occurrence, meaning that no symbol appears more than once in any program  $z \in \mathcal{L}$ . If there exist  $a \neq b$  such that*

$$\Delta_{\text{ctx}}(a, b; \eta^*) = 0,$$

*then there exists  $\tilde{\eta} \not\sim \eta^*$  such that*

$$Q_{\tilde{\eta}}^{(n)} = Q_{\eta^*}^{(n)} \quad \text{for every } n \geq 1.$$

So  $\Delta_{\text{ctx}} = 0$  is not merely a difficult regime: in single-occurrence languages it yields exact observational collapse.

Let

$$\bar{h}_X^2(\eta, \eta^*) := \frac{1}{m} \sum_{i=1}^m h^2(q_\eta(\cdot \mid X_i), q_{\eta^*}(\cdot \mid X_i))$$

denote the empirical conditional Hellinger distance in the reduced experiment.

**Theorem 3.3** (Predictive posterior contraction). *Assume the model and prior conditions of Appendix A, and suppose the design regularity regime*

$$me^{-ca, \delta^n} \rightarrow 0 \quad \text{for some fixed } 0 < \delta < 1.$$

*Then there exists  $M_{\text{pred}} \in (0, \infty)$  such that for every  $M \geq M_{\text{pred}}$ ,*

$$\Pi_m(\eta : \bar{h}_X(\eta, \eta^*) > M\delta_{m,n}) \rightarrow 0 \quad (8)$$

*in  $P_{\eta^*}$ -probability.*

The grouped predictive radius is therefore governed by the finite-dimensional program-law complexity  $K_{\text{eff}}$  and the ambient module dimension  $rd^2$ .

**Proposition 3.4** (Oracle latent-program recovery). *Assume  $\gamma_{\text{comp}}(\eta^*) > 0$ . Then on the high-probability design event  $\mathcal{G}_{m,n}(\delta)$ , for every  $i \in [m]$ ,*

$$\begin{aligned} \mathbb{E}_{\eta^*} \left[ 1 - \Pi_{\eta^*} \left( Z_i = Z_i^* \mid X_i, \widehat{W}_i \right) \mid X_i, Z_i^* \right] \\ \leq C_{\text{lat}} \exp \left( - \frac{(1-\delta)n \gamma_{\text{comp}}(\eta^*)^2}{8\sigma^2} \right), \end{aligned} \quad (9)$$

and the same bound holds after averaging over  $i = 1, \dots, m$ .

This isolates the within-task information scale for latent-program recovery:  $n\gamma_{\text{comp}}(\eta^*)^2$ .

### 3.2. Structural localization and minimax lower bounds

Predictive contraction does not by itself localize the latent grammar. The missing ingredient is a local inverse inequality from predictive distance back to quotient parameter distance.

**Theorem 3.5** (Structural contraction from a local inverse inequality). *Assume the design regime of Theorem 3.3. Suppose there exist constants*

$$\kappa \geq 1, \quad \rho_\star > 0, \quad c_\star > 0, \quad \underline{h}_\star > 0$$

such that, on  $\mathcal{G}_{m,n}(\delta)$ ,

$$d_q(\eta, \eta^*) \leq \rho_\star \implies \bar{h}_X^2(\eta, \eta^*) \geq c_\star d_q(\eta, \eta^*)^{2\kappa}, \quad (10)$$

$$d_q(\eta, \eta^*) \geq \rho_\star \implies \bar{h}_X(\eta, \eta^*) \geq \underline{h}_\star. \quad (11)$$

Then there exists  $M_{\text{str}} \in (0, \infty)$  such that for every  $M \geq M_{\text{str}}$ ,

$$\Pi_m \left( \eta : d_q(\eta, \eta^*) > M\delta_{m,n}^{1/\kappa} \right) \longrightarrow 0 \quad (12)$$

in  $P_{\eta^*}$ -probability.

**Corollary 3.6** (Regular structural contraction). *Suppose*

$$\lambda_{\text{lin}}(\eta^*) > 0, \quad \gamma_{\text{comp}}(\eta^*) > 0,$$

and assume the global separation condition Equation (11). Then the local inverse exponent is  $\kappa = 1$ , and there exists  $M_{\text{reg}} \in (0, \infty)$  such that

$$\Pi_m(\eta : d_q(\eta, \eta^*) > M\delta_{m,n}) \longrightarrow 0 \quad (13)$$

for every  $M \geq M_{\text{reg}}$ .

**Proposition 3.7** (Explicit quadratic singularity). *Consider the duplicated family*

$$\eta_t = \left( \left( \frac{1}{2}, \frac{1}{2} \right), (1+t, 1-t) \right), \quad t \geq 0,$$

with truth  $t^* = 0$ . Then

$$d_q(\eta_t, \eta_0) = 2t, \quad \bar{h}_X^2(\eta_t, \eta_0) \asymp n^2 d_q(\eta_t, \eta_0)^4$$

locally around  $t = 0$ . Consequently, there exists  $M_{\text{sing}} \in (0, \infty)$  such that

$$\Pi_m \left( \eta_t : d_q(\eta_t, \eta_0) > M \frac{\delta_{m,n}^{1/2}}{\sqrt{n}} \right) \longrightarrow 0 \quad (14)$$

for every  $M \geq M_{\text{sing}}$ .

This is the sharp structural gap of the paper: predictive learning can remain regular while grammar recovery is singular.

For a parameter class  $\mathcal{H} \subseteq \Xi(B, \mathcal{L})$ , write

$$\mathfrak{R}_{m,n}(\mathcal{H}; d_q) := \inf_{\hat{\eta}} \sup_{\eta \in \mathcal{H}} \mathbb{E}_{\eta} [d_q(\hat{\eta}, \eta)].$$

**Theorem 3.8** (Minimax lower bounds). *The following lower bounds hold.*

1. **Grammar-law complexity.** *There exists a finite submodel  $\mathcal{H}_{\text{gram}} \subseteq \Xi(B, \mathcal{L})$  with a fixed module library such that*

$$\mathfrak{R}_{m,n}(\mathcal{H}_{\text{gram}}; d_q) \gtrsim \sqrt{\frac{K_{\text{eff}}}{m}}. \quad (15)$$

2. **Anchored modules.** *If  $\mathcal{L}$  contains an  $s$ -anchor set, then there exists a finite anchored submodel  $\mathcal{H}_{\text{anc}}$  such that*

$$\mathfrak{R}_{m,n}(\mathcal{H}_{\text{anc}}; d_q) \gtrsim \sigma d \sqrt{\frac{s}{mn}}. \quad (16)$$

3. **Zero observability.** *If there exist  $\eta_0, \eta_1$  with*

$$\begin{aligned} d_q(\eta_0, \eta_1) = \Delta_0 > 0, \\ Q_{\eta_0}^{(n)} = Q_{\eta_1}^{(n)} \quad \text{for all } n, \end{aligned} \quad (17)$$

then

$$\mathfrak{R}_{m,n}(\{\eta_0, \eta_1\}; d_q) \geq \frac{\Delta_0}{4}. \quad (18)$$

4. **Singular duplicated family.** *In the duplicated family of Theorem 3.7,*

$$\begin{aligned} \mathfrak{R}_{m,n}(\{\eta_0, \eta_{t_{m,n}}\}; d_q) \gtrsim \frac{\sigma}{m^{1/4} \sqrt{n}}, \\ t_{m,n} \asymp \frac{\sigma}{m^{1/4} \sqrt{n}}. \end{aligned} \quad (19)$$

The four parts of Theorem 3.8 expose the statistical phase diagram of the problem. The grammar law contributes an unavoidable  $\sqrt{K_{\text{eff}}/m}$  term. Explicit anchor tasks recover the classical  $\sqrt{1/(mn)}$  shared-parameter scaling. Zero contextual observability creates an exact non-identifiability floor. And overfitted duplicated-module families are singular, with the sharp  $m^{-1/4}n^{-1/2}$  structural rate. These are precisely the hardness axes used to build the benchmark in Section 5.

## 4. Proof Architecture

The proof architecture has two logically distinct layers:

$$\begin{aligned} (X_i, Y_i)_{i=1}^m &\equiv (X_i, \widehat{W}_i)_{i=1}^m \\ &\implies \bar{h}_X(\eta, \eta^*) \lesssim \delta_{m,n} \\ &\implies d_q(\eta, \eta^*) \lesssim \delta_{m,n}^{1/\kappa}. \end{aligned}$$

The first implication is an exact reduction to a grouped non-i.i.d. mixture experiment; the second is a local inverse problem whose exponent is regular or singular. This is the same conceptual split that appears in singular statistical models and weakly identifiable mixtures, but here the latent object is a *task-family grammar* rather than an ordinary mixture component.

**Reduction and identifiability.** The reduction theorem shows that, conditional on the design  $X_i$  and latent program  $Z_i = z$ , the least-squares statistic  $\widehat{W}_i = Y_i X_i^\top (X_i X_i^\top)^{-1}$  is matrix normal around the composed operator  $A_z$ , while the residual is conditionally independent of  $\eta$ . Thus the posterior given grouped data is exactly the posterior given  $(X_i, \widehat{W}_i)_{i=1}^m$ . Identifiability then has two layers: the grouped law first determines the induced mixing measure on composed operators by common-covariance kernel injectivity, and  $\lambda_{\text{lin}}(\eta^*) > 0$  with  $\gamma_{\text{comp}}(\eta^*) > 0$  upgrades that statement to local quotient identifiability of the grammar itself.

**Predictive contraction and decoding.** Conditional on  $X_{1:m}$ , the reduced observations are independent but non-identically distributed, so the relevant predictive metric is the empirical conditional Hellinger distance  $\bar{h}_X$ . The predictive theorem is proved by verifying the Ghosal–van der Vaart conditions in this reduced experiment: exponentially consistent tests for  $\bar{h}_X$ -separated alternatives, entropy bounds at radius  $\delta_{m,n}$ , and local prior thickness of Kullback–Leibler neighborhoods. Once the true task-family parameter is fixed, latent-program recovery reduces to responsibilities in a finite Gaussian mixture, and the decoding error decays at the within-task information scale  $n\gamma_{\text{comp}}(\eta^*)^2$ .

**Inverse geometry and lower bounds.** Structural contraction is an inverse problem. In the regular regime, deconvolved Fourier witnesses produce a local inverse inequality of the form  $\bar{h}_X^2(\eta, \eta^*) \gtrsim d_q(\eta, \eta^*)^2$ . In the duplicated family, by contrast, the first-order term cancels exactly and the reduced  $\chi^2$ -divergence admits a closed-form quartic expansion, yielding  $\bar{h}_X^2(\eta_t, \eta_0) \asymp n^2 d_q(\eta_t, \eta_0)^4$  and hence  $\kappa = 2$ . The lower bounds then follow from explicit subexperiments: Fano packings for grammar-law complexity, oracle anchored submodels for exposed modules, an exact Le Cam two-point argument for zero observability, and a product- $\chi^2$  calculation for the duplicated singular family.

## 5. Theory-Native Benchmark and Experiments

Our evaluation philosophy is *theorem-native*: the benchmark is not an external testbed to which theory is retrofitted, but a direct instantiation of the stochastic task-grammar model analyzed above. Each benchmark cell is generated from a known latent task family and is released with exact hardness labels

$$(K_{\text{eff}}, \Delta_{\text{ctx}}, \lambda_{\text{lin}}, \gamma_{\text{comp}}, \kappa, s_{\text{anc}}),$$

computed from the definitions in Section 3 rather than inferred post hoc. In particular, the benchmark generator is designed so that the hardness axes correspond exactly to the statistical phase structure isolated by the theorems.

All experiments are run in the reduced experiment  $(X_i, \widehat{W}_i)_{i=1}^m$ , which is Blackwell equivalent to the full grouped-data experiment by Theorem 2.1. This keeps the empirical section aligned with the proved model and removes optimization confounders unrelated to the statistical questions studied here. Posterior inference is exact by low-dimensional grid integration in the observability and singular suites, and MAP plus exact-Hessian Laplace approximation in the theorem-matched regular and anchored suites. The primary metrics are  $R_{\text{pred},0.9}$ ,  $R_{\text{str},0.9}$ ,  $R_{\text{anc},0.9}$ , and  $M_Z := m^{-1} \sum_{i=1}^m \Pi_m(Z_i = Z_i^*)$ ; curves show medians over 32 replications.

To conserve main-paper space, we place the regular-family sanity checks in Appendix I. Those figures verify that, in the regular regime, normalization by the theoretical predictive and structural scales substantially stabilizes the posterior radii, and that the  $q = 4$  family exhibits the predicted local inverse law  $\bar{h}_X(\eta, \eta^*)^2 \asymp d_q(\eta, \eta^*)^2$ . The main text instead concentrates on the three benchmark regimes that most sharply diagnose the new hardness parameters: observability, singularity, and anchor exposure.

**Observability phase transition.** We use

$$\mathcal{L}_{\text{obs}} = \{(1, 3), (2, 3)\}, \quad \pi^* = \left(\frac{1}{2}, \frac{1}{2}\right),$$

with

$$A_1^* = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \quad A_2^* = \begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix}, \quad A_3^*(\rho) = \begin{bmatrix} 1 & 0 \\ 0 & \rho \end{bmatrix}.$$

Then  $\Delta_{\text{ctx}}(\eta_\rho^*) = \rho^2/2$  and  $\gamma_{\text{comp}}(\eta_\rho^*) = \rho$ . Figure 1 shows the predicted collapse of the latent-program error  $1 - M_Z$  against the theorem scale  $n\Delta_{\text{ctx}}$ . At  $\rho = 0$ , both  $R_{\text{str},0.9}$  and  $1 - M_Z$  remain bounded away from zero as  $n$  increases, matching Theorems 3.2 and 3.8.

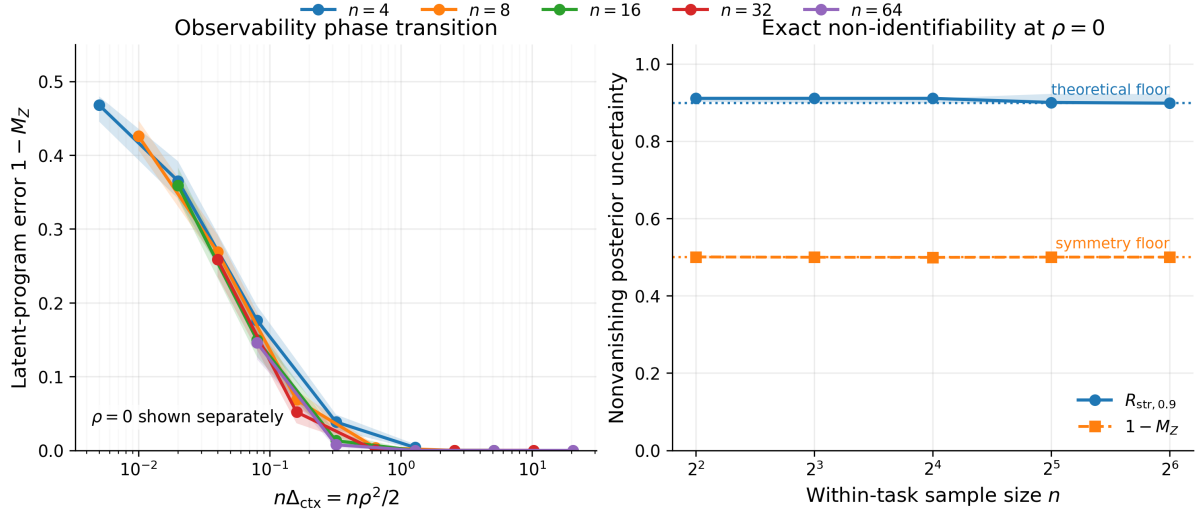


Figure 1. **Observability phase transition.** Left: latent-program error collapses when plotted against the theorem scale  $n\Delta_{\text{ctx}} = n\rho^2/2$ , confirming that the within-task information budget is the relevant observable quantity. Right: at  $\rho = 0$ , both  $R_{\text{str},0.9}$  and  $1 - M_Z$  remain bounded away from zero, confirming exact non-identifiability. Curves show medians; bands are 10th–90th percentiles.

**Singularity gap: prediction is regular, structure is singular.** We next evaluate the duplicated-module family

$$\mathcal{L}_{\text{sing}} = \{(1), (2)\}, \quad \pi^* = \left(\frac{1}{2}, \frac{1}{2}\right), \\ A_1(t) = 1 + t, \quad A_2(t) = 1 - t$$

with truth  $t^* = 0$ . This is the explicit  $\kappa = 2$  family of Theorem 3.7. The theory predicts

$$R_{\text{pred},0.9} = O(\delta_{m,n}), \quad R_{\text{str},0.9} = O\left(\frac{\delta_{m,n}^{1/2}}{\sqrt{n}}\right).$$

Figure 2 confirms both claims and directly validates the quartic inverse law  $\bar{h}_X(\eta_t, \eta_0)^2 \asymp n^2 d_q(\eta_t, \eta_0)^4$ .

**Anchors as benchmark design interventions.** Finally, consider

$$\mathcal{L}_s^{\text{anc}} = \{(1), \dots, (s)\} \cup \{(s+1), (s+1, s+2)\},$$

where the first  $s$  programs expose anchor modules in isolation, and the matched no-anchor family

$$\mathcal{L}_s^{\text{noanc}} = \{(1, s+1), \dots, (s, s+1)\} \\ \cup \{(s+1), (s+1, s+2)\},$$

where the same modules appear only through a shared context. By Theorem 3.8, the natural normalized anchor error is

$$\frac{R_{\text{anc},0.9}\sqrt{mn}}{d s^{3/2}}.$$

The left panel of Figure 3 shows that this quantity stays order one across the anchor-count grid. The right panel

isolates the benchmark-design effect: explicit anchor tasks materially reduce anchor-module posterior radius relative to matched no-anchor cells.

Taken together, the main-text experiments validate three theorem-guided hardness mechanisms: observability controls the within-task phase transition for latent-program recovery, singularity separates predictive and structural learning, and anchor exposure changes which latent modules are directly identifiable from benchmark cells. The benchmark is therefore not merely compatible with the theory; it is organized by the theory.

These three suites also isolate genuinely different statistical resources. In the observability family the grammar law is fixed and only the within-task information budget  $n\Delta_{\text{ctx}}$  is varied, so the phase transition cannot be mistaken for poor cross-task coverage. In the singular family the grouped predictive experiment is regular while the local inverse map is quartic, so larger  $m$  improves predictive fit faster than it resolves latent structure. In the anchor family the ambient noise level and module class are held fixed while only exposure pattern changes. That separation is why the plots should be read as benchmark-design interventions on identifiability, not merely as descriptive summaries of synthetic data.

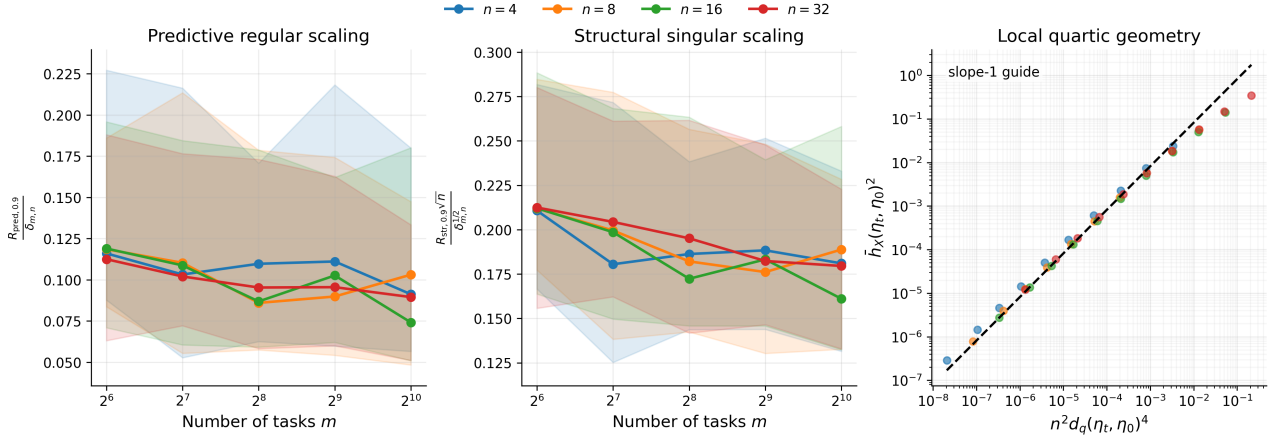


Figure 2. Singularity gap in the duplicated-module family. Left: after normalization by the regular predictive scale  $\delta_{m,n}$ , the predictive posterior radius  $R_{\text{pred},0.9}$  is approximately stable across  $m$ , consistent with regular predictive learning. Middle: after normalization by the singular structural scale  $\delta_{m,n}^{1/2}/\sqrt{n}$ , the structural posterior radius  $R_{\text{str},0.9}$  is approximately stable across  $m$ , confirming slower structural contraction. Right: direct local-geometry validation of the quartic law  $\bar{h}_X^2 \propto n^2 d_q^4$  underlying  $\kappa = 2$ . Curves show medians; bands are 10th–90th percentiles.

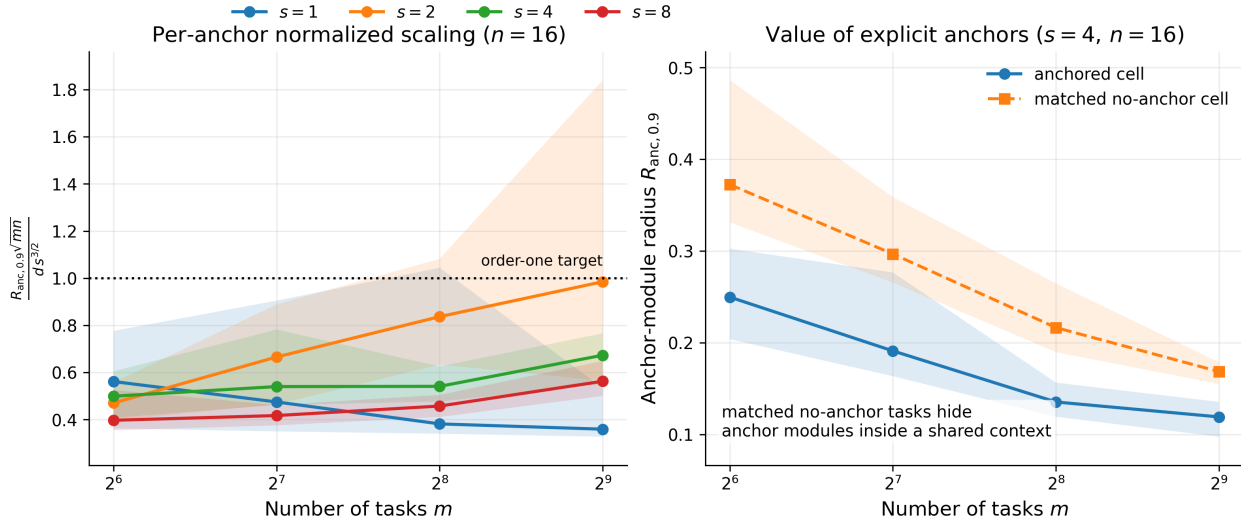


Figure 3. Anchored-module scaling and the value of explicit anchors. Left: after per-anchor normalization, the anchor-module radius remains order one across the anchor-count grid, consistent with the predicted  $d s^{3/2}/\sqrt{mn}$  scaling for the summed anchor-module error. Right: benchmark cells with explicit anchor tasks yield substantially smaller anchor-module posterior radius than matched no-anchor cells, in which anchor modules appear only inside a shared context. Curves show medians; bands are interquartile ranges.

## 6. Discussion

The main conceptual shift in this paper is to treat the *task-generating family*, rather than the individual task, as the statistical object of interest. Once tasks are modeled as draws from a latent compositional grammar, transfer across related tasks, uncertainty over latent task identity, module reuse, and the effect of benchmark design on identifiability become part of one posterior problem. In that formulation, effective grammar complexity, contextual observability, and singularity are not metaphors; they are explicit hardness parameters that govern upper and lower bounds.

A second message is that *prediction and structure are not the same learning problem*. In regular families, predictive contraction and structural localization occur at the same order. In singular families, the inverse map from predictive laws back to latent task structure flattens, so structural uncertainty can decay strictly more slowly than predictive uncertainty. For evaluation, this means that strong held-out performance on related tasks does not by itself certify rapid recovery of the reusable compositional structure.

A third message is methodological: benchmark design should be viewed as a *statistical intervention*. In our analysis, observability, anchor exposure, and component separation are not cosmetic properties of a benchmark suite; they determine identifiability and rates. This viewpoint is complementary to recent work on benchmark compression and score prediction (Procaccia et al., 2025; Zhao et al., 2024): rather than selecting representative metrics or imputing missing scores after a benchmark is fixed, we expose ex ante hardness knobs that can be dialed deliberately when constructing benchmark cells.

Our results are proved for finite-language stochastic linear grammars, and the experiments validate the theory inside theorem-matched synthetic subfamilies. We therefore do not claim that present-day foundation models are exactly described by the model class studied here. The point is more structural: if one wants guarantees for transfer across compositional task families, then observability, inverse geometry, and benchmark-cell design must enter the theory explicitly.

From that perspective, the benchmark generator should be read as a controlled diagnostic instrument. Grammar-law complexity determines how many distinct tasks are needed before the latent program law can be learned. Observability and component separation determine whether within-task data are even informative enough to decode latent programs. Singularity determines whether predictive success can be converted into structural certainty. Anchor exposure determines which parts of the library are directly visible to the learner. Separating these mechanisms matters in practice: they correspond to different failure modes, different remedies, and different claims that an evaluation suite can responsibly support.

The most important next step is to extend this theorem-benchmark cycle beyond linear grammars without losing the quotient and inverse-geometry viewpoint. Natural directions include bounded-arity tree grammars, nonlinear module classes, controlled model misspecification, and larger benchmark suites whose cells are labeled ex ante by statistical hardness rather than only scored ex post by predictive accuracy. Even when the present model is only an approximation, the distinction between predictive and structural uncertainty, and the role of benchmark design in mediating that distinction, should remain central.

The rates also separate cross-task and within-task sample budgets. The term  $K_{\text{eff}} \log m/m$  is a grammar-law burden that only additional tasks can reduce, whereas observability and singularity operate through the within-task experiment and the local inverse geometry. A suite can therefore be rich in examples per task yet underpowered for recovering the law over programs, or conversely span many tasks while still revealing too little within-task information to decode latent structure. For theorem-native evaluation,  $m$  and  $n$  should be treated as distinct resources rather than collapsed into a single total sample size.

One practical implication is that theorem-native benchmark cells can disaggregate failure modes that are otherwise collapsed into a single held-out score. Poor performance may reflect insufficient cross-task coverage of the grammar law, insufficient within-task evidence to decode latent programs, hidden singularity in the inverse map from predictive behavior back to structure, or the simple absence of anchor exposure for the modules one hopes to identify. Those

cases call for different responses. More tasks help the first, redesigned cells help the second and fourth, and the third requires acknowledging persistent structural uncertainty even when predictive fit is strong.

Finally, the quotient formulation clarifies what successful recovery can mean. Even in favorable regular regimes the target is an equivalence class under  $\text{Aut}(\mathcal{L})$ , not a preferred symbolic labeling of modules. Conversely, when grouped laws coincide for reasons not induced by an automorphism, the ambiguity is substantive and should be reported as non-identifiability. The  $\rho = 0$  observability cell is informative precisely because it separates unavoidable statistical collapse from mere relabeling symmetry.

More broadly, we view the present linear setting as an existence proof that benchmark construction can itself be placed on a statistical footing. The lesson is not that realistic foundation-model pipelines are linear, but that benchmark suites should ideally be accompanied by formal statements about what their cells reveal and what they cannot reveal. In that sense, hardness labels play the role of experimental metadata: they make evaluation results easier to interpret, easier to compare across cells, and harder to overstate.

## Impact Statement

This paper develops a statistical theory of compositional task families and uses that theory to construct theorem-native benchmark cells with explicit hardness labels. The main positive impact is methodological. Our results suggest that benchmark design should expose the structural quantities that determine what can and cannot be learned—for example observability, anchor exposure, and singularity—rather than relying only on aggregate held-out accuracy. In settings where one wants guarantees for transfer across related tasks, this can lead to more informative evaluation suites, finer-grained failure analysis, and better calibrated uncertainty about latent reusable structure. Because the benchmark families studied here are synthetic and theorem-matched, the work does not directly increase model capabilities or provide a recipe for scaling a deployed system.

The main risk is *over-interpretation*. Our guarantees are proved for finite-language stochastic linear grammars and experimentally validated only within theorem-matched synthetic regimes. If the results were treated as applying unchanged to realistic foundation-model pipelines, they could create false confidence about identifiability, transfer, or uncertainty decomposition in settings that violate the model assumptions. A related risk is that theorem-native benchmark generators, if used carelessly, could themselves become targets for narrow benchmark-specific optimization rather than broader scientific understanding. We therefore view the intended use of this work as *diagnostic* rather than *certificatory*: hardness labels should guide evaluation design and hypothesis testing, not substitute for robustness checks, external validation, or real-world safety assessment.

Finally, the paper may help reduce a common failure mode in evaluation: conflating predictive success with structural understanding. Our singular-family analysis shows that a model can achieve regular predictive contraction while remaining substantially uncertain about the underlying compositional structure. Making that distinction explicit is, in our view, socially beneficial, because it pushes benchmark practice away from overclaiming what a finite suite of successful evaluations actually certifies.

## References

- Abedsoltan, A., Zhang, H., Wen, K., Lin, H., Zhang, J., and Belkin, M. Task generalization with autoregressive compositional structure: Can learning from  $D$  tasks generalize to  $D^T$  tasks? In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 154–173. PMLR, 2025. URL <https://proceedings.mlr.press/v267/abedsoltan25a.html>.
- Birgé, L. Robust tests for model selection. In Banerjee, M., Bunea, F., Huang, J., Koltchinskii, V., and Maathuis, M. H. (eds.), *From Probability to Statistics and Back: High-Dimensional Models and Processes*, volume 9 of *IMS Collections*, pp. 47–64. Institute of Mathematical Statistics, 2013. doi: 10.1214/12-IMSCOLL905.
- Chang, H., Park, J., Cho, H., Yang, S., Ko, M., Hwang, H., Won, S., Lee, D., Ahn, Y., and Seo, M. The coverage principle: A framework for understanding compositional generalization. *arXiv preprint arXiv:2505.20278*, 2025. doi: 10.48550/arXiv.2505.20278.
- Ghosal, S. and van der Vaart, A. W. Convergence rates of posterior distributions for non-i.i.d. observations. *The Annals of Statistics*, 35(1):192–223, 2007. doi: 10.1214/009053606000001172.
- Ghosal, S., Ghosh, J. K., and van der Vaart, A. W. Convergence rates of posterior distributions. *The Annals of Statistics*, 28(2):500–531, 2000. doi: 10.1214/aos/1016218228.
- Ho, N. and Nguyen, X. Convergence rates of parameter estimation for some weakly identifiable finite mixtures. *The Annals of Statistics*, 44(6):2726–2755, 2016. doi: 10.1214/16-AOS1444.
- Ho, N. and Nguyen, X. Singularity structures and impacts on parameter estimation in finite mixtures of distributions. *SIAM Journal on Mathematics of Data Science*, 1(4):730–758, 2019. doi: 10.1137/18M122947X.
- Procaccia, A., Schiffer, B., Wang, S., and Zhang, S. Metriocracy: Representative metrics for lite benchmarks. *arXiv preprint arXiv:2506.09813*, 2025. doi: 10.48550/arXiv.2506.09813.
- Schug, S., Kobayashi, S., Akram, Y., Wołczyk, M., Proca, A., von Oswald, J., Pascanu, R., Sacramento, J., and Steger, A. Discovering modular solutions that generalize compositionally. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=H98CVcX1eh>.
- Teicher, H. Identifiability of finite mixtures. *The Annals of Mathematical Statistics*, 34(4):1265–1269, 1963. doi: 10.1214/aoms/1177703862.
- Tsybakov, A. B. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer, 2009. doi: 10.1007/b13794.
- Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018. doi: 10.1017/9781108231596.
- Watanabe, S. *Algebraic Geometry and Statistical Learning Theory*. Cambridge University Press, 2009. doi: 10.1017/CBO9780511800474.
- Yakowitz, S. J. and Spragins, J. D. On the identifiability of finite mixtures. *The Annals of Mathematical Statistics*, 39(1):209–214, 1968. doi: 10.1214/aoms/1177698520.
- Zhao, Q., Xu, M., Gupta, K., Asthana, A., Zheng, L., and Gould, S. Can we predict performance of large models across vision-language tasks? *arXiv preprint arXiv:2410.10112*, 2024. doi: 10.48550/arXiv.2410.10112.

## A. Additional Notation, Assumptions, and Quotient Geometry

This appendix freezes the exact parameter space, group action, and hardness functionals used throughout the paper. No later section changes these definitions.

### A.1. Parameter space and basic notation

Fix an admissible language  $\mathcal{L} \subseteq \mathcal{Z}_{L,r}$ , where

$$\mathcal{Z}_{L,r} = \bigcup_{t=1}^L [r]^t, \quad [r] = \{1, \dots, r\}, \quad [r]^0 = \{\emptyset\}.$$

For  $z = (z_1, \dots, z_t) \in \mathcal{Z}_{L,r}$ , let  $|z| = t$ . If  $u = (u_1, \dots, u_s) \in [r]^s$  and  $v = (v_1, \dots, v_q) \in [r]^q$ , we write  $uav$  for the concatenated string

$$(u_1, \dots, u_s, a, v_1, \dots, v_q) \in [r]^{s+1+q}.$$

The grammar parameter space is

$$\mathfrak{P}_{\mathcal{L}}^{\circ} = \left\{ \pi = (\pi_z)_{z \in \mathcal{L}} \in (0, 1)^{\mathcal{L}} : \sum_{z \in \mathcal{L}} \pi_z = 1 \right\},$$

and the bounded module space is

$$\mathbb{A}_B = \{A \in \mathbb{R}^{d \times d} : \|A\|_{\text{op}} \leq B\}.$$

The full parameter space is

$$\Xi(B, \mathcal{L}) = \mathfrak{P}_{\mathcal{L}}^{\circ} \times \mathbb{A}_B^r.$$

If  $\Theta = (A_1, \dots, A_r) \in \mathbb{A}_B^r$ , then for any nonempty string  $z = (z_1, \dots, z_t)$  we define

$$A_z = A_{z_t} \cdots A_{z_1}, \quad A_{\emptyset} = I_d.$$

For any  $W \in \mathbb{R}^{d \times d}$ ,  $P_W^{(n)}$  denotes the law of one task dataset  $D = (X, Y)$  under

$$x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d), \quad y_j = Wx_j + \varepsilon_j, \quad \varepsilon_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2 I_d),$$

with the Gaussian design and noise independent.

For dominated probability laws  $P$  and  $Q$  with densities  $p$  and  $q$ , respectively, we write

$$h^2(P, Q) = \frac{1}{2} \int (\sqrt{p} - \sqrt{q})^2 d\mu, \quad \text{KL}(P\|Q) = \int p \log \frac{p}{q} d\mu,$$

where  $\mu$  is any common dominating measure.

### A.2. Standing model and prior conditions

**Assumption A.1** (Sampling model and truth). Throughout, the admissible language  $\mathcal{L} \subseteq \mathcal{Z}_{L,r}$ , the dimensions  $(d, r, L)$ , the operator bound  $B$ , and the noise level  $\sigma > 0$  are fixed. The true parameter is  $\eta^* = (\pi^*, \Theta^*) \in \Xi(B, \mathcal{L})$ , with  $\pi_z^* \geq \pi_{\min} > 0$  for every  $z \in \mathcal{L}$ . For each  $z \in \mathcal{L}$ , let  $A_z^*$  denote the composed operator induced by  $\Theta^*$ . For each task  $i \in [m]$ , the latent program satisfies  $Z_i \sim \pi^*$ , the design vectors satisfy  $x_{i1}, \dots, x_{in} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$  with  $n \geq d + 1$ , and, conditional on  $Z_i = z$ , the responses obey

$$y_{ij} = A_z^* x_{ij} + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2 I_d),$$

with the Gaussian designs, latent programs, and noises mutually independent across tasks.

**Assumption A.2** (Prior regularity). The prior factors as  $\Pi = \Pi_{\pi} \otimes \Pi_{\Theta}$  on  $\Xi(B, \mathcal{L})$ . The grammar-law marginal  $\Pi_{\pi}$  admits a density with respect to the intrinsic Lebesgue measure on  $\mathfrak{P}_{\mathcal{L}}^{\circ}$  that is continuous and strictly positive in a neighborhood of  $\pi^*$ . The library marginal  $\Pi_{\Theta} = \bigotimes_{a=1}^r \Pi_a$  admits a density on  $\mathbb{A}_B^r$  that is continuous and strictly positive in a neighborhood of  $\Theta^*$ .

### A.3. Automorphism group and quotient metric

The correct symmetry group is the automorphism group of the admissible language.

**Definition A.3** (Automorphism group). The automorphism group of  $\mathcal{L}$  is

$$\text{Aut}(\mathcal{L}) = \{\sigma \in S_r : \sigma \cdot \mathcal{L} = \mathcal{L}\},$$

where  $\sigma \cdot z = (\sigma(z_1), \dots, \sigma(z_t))$  for every string  $z = (z_1, \dots, z_t)$ .

For  $\sigma \in \text{Aut}(\mathcal{L})$  and  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$ , define

$$(\sigma \cdot \pi)_u := \pi_{\sigma^{-1}(u)}, \quad (\sigma \cdot \Theta)_a := A_{\sigma^{-1}(a)}, \quad \sigma \cdot \eta := (\sigma \cdot \pi, \sigma \cdot \Theta).$$

The base distance on  $\Xi(B, \mathcal{L})$  is

$$d_{\text{par}}((\pi, \Theta), (\pi', \Theta')) = \|\pi - \pi'\|_1 + \sum_{a=1}^r \|A_a - A'_a\|_F.$$

The quotient distance is

$$d_q(\eta, \eta') = \min_{\sigma \in \text{Aut}(\mathcal{L})} d_{\text{par}}(\eta, \sigma \cdot \eta').$$

The next lemma shows that the automorphism action preserves the grouped law exactly.

**Lemma A.4** (Permutation invariance of the grouped law). For every  $\eta \in \Xi(B, \mathcal{L})$ , every  $\sigma \in \text{Aut}(\mathcal{L})$ , and every  $n \geq 1$ ,

$$Q_{\sigma \cdot \eta}^{(n)} = Q_{\eta}^{(n)}.$$

*Proof.* Write  $\eta = (\pi, \Theta)$ . For any  $z = (z_1, \dots, z_t) \in \mathcal{L}$ ,

$$A_{\sigma \cdot z}(\sigma \cdot \Theta) = (\sigma \cdot \Theta)_{\sigma(z_t)} \cdots (\sigma \cdot \Theta)_{\sigma(z_1)} = A_{z_t} \cdots A_{z_1} = A_z(\Theta).$$

Since  $\sigma \in \text{Aut}(\mathcal{L})$ , the map  $z \mapsto \sigma \cdot z$  is a bijection of  $\mathcal{L}$ . Therefore,

$$Q_{\sigma \cdot \eta}^{(n)} = \sum_{u \in \mathcal{L}} (\sigma \cdot \pi)_u P_{A_u(\sigma \cdot \Theta)}^{(n)} = \sum_{z \in \mathcal{L}} \pi_z P_{A_{\sigma \cdot z}(\sigma \cdot \Theta)}^{(n)} = \sum_{z \in \mathcal{L}} \pi_z P_{A_z(\Theta)}^{(n)} = Q_{\eta}^{(n)}.$$

□

**Lemma A.5** (The quotient distance is well defined). The function  $d_q$  induces a metric on the quotient space  $\Xi(B, \mathcal{L}) / \text{Aut}(\mathcal{L})$ .

*Proof.* We first note that the base distance is invariant under the simultaneous group action:

$$d_{\text{par}}(\sigma \cdot \eta, \sigma \cdot \eta') = d_{\text{par}}(\eta, \eta') \quad \text{for all } \sigma \in \text{Aut}(\mathcal{L}).$$

Indeed, permutation of the coordinates of  $\pi$  preserves the  $\ell_1$ -norm, and permutation of the module labels preserves  $\sum_{a=1}^r \|A_a - A'_a\|_F$ .

Nonnegativity of  $d_q$  is immediate. If  $d_q(\eta, \eta') = 0$ , then because  $\text{Aut}(\mathcal{L})$  is finite, the minimum is attained by some  $\sigma \in \text{Aut}(\mathcal{L})$ , and

$$0 = d_{\text{par}}(\eta, \sigma \cdot \eta')$$

implies  $\eta = \sigma \cdot \eta'$ . Hence  $\eta$  and  $\eta'$  represent the same quotient class.

For symmetry, using invariance of  $d_{\text{par}}$  and the fact that  $\sigma^{-1} \in \text{Aut}(\mathcal{L})$  whenever  $\sigma \in \text{Aut}(\mathcal{L})$ ,

$$d_q(\eta, \eta') = \min_{\sigma} d_{\text{par}}(\eta, \sigma \cdot \eta') = \min_{\sigma} d_{\text{par}}(\sigma^{-1} \cdot \eta, \eta') = d_q(\eta', \eta).$$

For the triangle inequality, let  $\eta, \eta', \eta'' \in \Xi(B, \mathcal{L})$ , and fix  $\sigma, \tau \in \text{Aut}(\mathcal{L})$ . Then

$$d_{\text{par}}(\eta, \sigma\tau \cdot \eta'') \leq d_{\text{par}}(\eta, \sigma \cdot \eta') + d_{\text{par}}(\sigma \cdot \eta', \sigma\tau \cdot \eta'') = d_{\text{par}}(\eta, \sigma \cdot \eta') + d_{\text{par}}(\eta', \tau \cdot \eta'').$$

Taking the minimum first over  $\sigma$ , then over  $\tau$ , yields

$$d_q(\eta, \eta'') \leq d_q(\eta, \eta') + d_q(\eta', \eta'').$$

Hence  $d_q$  is a metric on the quotient space. □

#### A.4. Context distribution and contextual observability

The contextual observability functional is defined using the distribution of one-hole contexts induced by the latent task-family law.

**Definition A.6** (One-hole contexts). The set of one-hole contexts associated with  $\mathcal{L}$  is

$$\mathcal{C}_{\mathcal{L}} = \left\{ (u, v) : u \in \bigcup_{s=0}^{L-1} [r]^s, v \in \bigcup_{q=0}^{L-1} [r]^q, \exists a \in [r] \text{ such that } uav \in \mathcal{L} \right\}.$$

For  $z = (z_1, \dots, z_t) \in \mathcal{L}$  and  $\ell \in \{1, \dots, t\}$ , define

$$\text{ctx}_{\ell}(z) = ((z_1, \dots, z_{\ell-1}), (z_{\ell+1}, \dots, z_t)) \in \mathcal{C}_{\mathcal{L}}.$$

**Lemma A.7** (The context law is a probability mass function). *For every  $\pi \in \mathfrak{P}_{\mathcal{L}}^{\circ}$ , the function*

$$\nu_{\pi}(u, v) = \sum_{z \in \mathcal{C}_{\mathcal{L}}} \pi_z \frac{1}{|z|} \sum_{\ell=1}^{|z|} \mathbf{1}\{\text{ctx}_{\ell}(z) = (u, v)\}, \quad (u, v) \in \mathcal{C}_{\mathcal{L}},$$

defines a probability mass function on  $\mathcal{C}_{\mathcal{L}}$ .

*Proof.* Nonnegativity is immediate. Summing over  $(u, v) \in \mathcal{C}_{\mathcal{L}}$  and exchanging the order of summation gives

$$\sum_{(u, v) \in \mathcal{C}_{\mathcal{L}}} \nu_{\pi}(u, v) = \sum_{z \in \mathcal{C}_{\mathcal{L}}} \pi_z \frac{1}{|z|} \sum_{\ell=1}^{|z|} \sum_{(u, v) \in \mathcal{C}_{\mathcal{L}}} \mathbf{1}\{\text{ctx}_{\ell}(z) = (u, v)\}.$$

For each fixed  $z$  and  $\ell$ , exactly one context is selected, so the innermost sum equals 1. Therefore,

$$\sum_{(u, v) \in \mathcal{C}_{\mathcal{L}}} \nu_{\pi}(u, v) = \sum_{z \in \mathcal{C}_{\mathcal{L}}} \pi_z \frac{1}{|z|} \sum_{\ell=1}^{|z|} 1 = \sum_{z \in \mathcal{C}_{\mathcal{L}}} \pi_z = 1.$$

□

For  $\eta = (\pi, \Theta)$ ,  $\Theta = (A_1, \dots, A_r)$ , and  $a \neq b$ , define

$$\Delta_{\text{ctx}}(a, b; \eta) = \sum_{(u, v) \in \mathcal{C}_{\mathcal{L}}} \nu_{\pi}(u, v) \|A_v(A_a - A_b)A_u\|_F^2, \quad \Delta_{\text{ctx}}(\eta) = \min_{a \neq b} \Delta_{\text{ctx}}(a, b; \eta).$$

**Lemma A.8** (Continuity of contextual observability). *For each pair  $a \neq b$ , the map*

$$\eta \mapsto \Delta_{\text{ctx}}(a, b; \eta)$$

is continuous on  $\Xi(B, \mathcal{L})$ . Consequently,  $\eta \mapsto \Delta_{\text{ctx}}(\eta)$  is continuous on  $\Xi(B, \mathcal{L})$ .

*Proof.* Because  $\mathcal{L}$  is finite,  $\mathcal{C}_{\mathcal{L}}$  is finite. For fixed  $(u, v)$ , the map

$$(\pi, \Theta) \mapsto \nu_{\pi}(u, v) \|A_v(A_a - A_b)A_u\|_F^2$$

is continuous:  $\nu_{\pi}(u, v)$  is linear in  $\pi$ , and  $A_v(A_a - A_b)A_u$  is a polynomial expression in the entries of the module matrices. Hence  $\Delta_{\text{ctx}}(a, b; \eta)$ , being a finite sum of continuous functions, is continuous. Since  $\Delta_{\text{ctx}}(\eta)$  is the minimum over finitely many continuous functions indexed by  $(a, b)$  with  $a \neq b$ , it is also continuous. □

**Lemma A.9** (Quotient invariance of contextual observability). *For every  $\eta \in \Xi(B, \mathcal{L})$  and every  $\sigma \in \text{Aut}(\mathcal{L})$ ,*

$$\Delta_{\text{ctx}}(\sigma \cdot \eta) = \Delta_{\text{ctx}}(\eta).$$

*Proof.* Fix distinct  $a, b \in [r]$ . For every  $u, v$ , the map  $(u, v) \mapsto (\sigma \cdot u, \sigma \cdot v)$  is a bijection of  $\mathcal{C}_{\mathcal{L}}$ , because  $\sigma \in \text{Aut}(\mathcal{L})$ . Moreover,

$$\nu_{\sigma \cdot \pi}(\sigma \cdot u, \sigma \cdot v) = \nu_{\pi}(u, v),$$

since relabeling the symbols of a latent program and then choosing a uniformly random position leaves the induced context weights unchanged after reindexing. Also,

$$A_{\sigma \cdot v}(\sigma \cdot \Theta) \left( (\sigma \cdot \Theta)_{\sigma(a)} - (\sigma \cdot \Theta)_{\sigma(b)} \right) A_{\sigma \cdot u}(\sigma \cdot \Theta) = A_v(\Theta) (A_a(\Theta) - A_b(\Theta)) A_u(\Theta).$$

Therefore,

$$\Delta_{\text{ctx}}(\sigma(a), \sigma(b); \sigma \cdot \eta) = \Delta_{\text{ctx}}(a, b; \eta).$$

Taking the minimum over all distinct pairs yields

$$\Delta_{\text{ctx}}(\sigma \cdot \eta) = \Delta_{\text{ctx}}(\eta).$$

□

*Remark A.10.* The quantities  $Q_{\eta}^{(n)}$ ,  $d_q$ , and  $\Delta_{\text{ctx}}(\eta)$  are therefore all well defined on quotient classes. Later, when the local inverse inequality is established, the singularity order  $\kappa(\eta^*)$  will also be a quotient-invariant object.

## B. Sufficient-Statistic Reduction

This appendix proves that, conditional on the random design, no inferential information about the latent task-family parameter  $\eta = (\pi, \Theta)$  is lost by replacing each raw task dataset  $D = (X, Y)$  with the task-level least-squares estimator  $\widehat{W}$ . The reduction is exact: under any prior on  $\eta$ , the posterior based on the full grouped data  $(X_i, Y_i)_{i=1}^m$  coincides almost surely with the posterior based on the reduced experiment  $(X_i, \widehat{W}_i)_{i=1}^m$ . Moreover, conditional on the designs, the reduced experiment is a heteroskedastic finite mixture of matrix-normal laws centered at the composed operators  $\{A_z\}_{z \in \mathcal{L}}$ .

### B.1. Task-level least-squares reduction

We first record the matrix-normal convention used throughout the appendix.

**Definition B.1** (Matrix-normal law). For integers  $p, q \geq 1$ , a random matrix  $M \in \mathbb{R}^{p \times q}$  is said to follow the matrix-normal law

$$M \sim \mathcal{MN}_{p,q}(M_0, U, V)$$

if

$$\text{vec}(M) \sim \mathcal{N}(\text{vec}(M_0), V \otimes U),$$

where  $M_0 \in \mathbb{R}^{p \times q}$ ,  $U \in \mathbb{R}^{p \times p}$  is positive semidefinite, and  $V \in \mathbb{R}^{q \times q}$  is positive semidefinite.

Fix one task  $D = (X, Y)$ , where

$$X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}, \quad Y = [y_1, \dots, y_n] \in \mathbb{R}^{d \times n}.$$

Recall from Theorem A.1 that  $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ , and  $n \geq d + 1$ .

Define the empirical Gram matrix

$$G_X := XX^{\top} \in \mathbb{R}^{d \times d}.$$

**Lemma B.2** (Full row rank of the Gaussian design). *Under Theorem A.1,  $G_X$  is positive definite almost surely.*

*Proof.* It suffices to show that  $X$  has row rank  $d$  almost surely. Consider the first  $d$  columns  $(x_1, \dots, x_d)$ . Since these are independent draws from an absolutely continuous distribution on  $\mathbb{R}^d$ , the  $d \times d$  matrix  $[x_1, \dots, x_d]$  has nonzero determinant almost surely. Hence  $X$  has row rank  $d$  almost surely, and therefore  $G_X = XX^{\top}$  is invertible almost surely. □

In view of Theorem B.2, all identities below hold on an event of probability one. We suppress this null set from the notation.

Define the task-level least-squares estimator

$$\widehat{W} := YX^\top G_X^{-1} \in \mathbb{R}^{d \times d}, \quad (20)$$

the design projection

$$P_X := X^\top G_X^{-1} X \in \mathbb{R}^{n \times n}, \quad (21)$$

and the residual matrix

$$R := Y(I_n - P_X) \in \mathbb{R}^{d \times n}. \quad (22)$$

**Proposition B.3** (Orthogonal decomposition of the task likelihood). *For every  $X$  with  $G_X$  invertible and every  $Y \in \mathbb{R}^{d \times n}$ , the following identities hold:*

1.  $Y = \widehat{W}X + R;$

2.  $RX^\top = 0;$

3. for every  $W \in \mathbb{R}^{d \times d},$

$$\|Y - WX\|_F^2 = \|(\widehat{W} - W)X\|_F^2 + \|R\|_F^2 = \text{tr}((\widehat{W} - W)G_X(\widehat{W} - W)^\top) + \|R\|_F^2. \quad (23)$$

*Proof.* By definition,

$$\widehat{W}X = YX^\top G_X^{-1} X = YP_X.$$

Hence

$$\widehat{W}X + R = YP_X + Y(I_n - P_X) = Y.$$

This proves part (1).

For part (2), using  $P_X X^\top = X^\top G_X^{-1} X X^\top = X^\top,$

$$RX^\top = Y(I_n - P_X)X^\top = Y(X^\top - P_X X^\top) = 0.$$

For part (3), fix  $W \in \mathbb{R}^{d \times d}.$  By part (1),

$$Y - WX = (\widehat{W} - W)X + R.$$

Using part (2),

$$\langle (\widehat{W} - W)X, R \rangle_F = \text{tr}((\widehat{W} - W)X R^\top) = \text{tr}(R X^\top (\widehat{W} - W)^\top) = 0.$$

Therefore

$$\|Y - WX\|_F^2 = \|(\widehat{W} - W)X\|_F^2 + \|R\|_F^2.$$

Finally,

$$\|(\widehat{W} - W)X\|_F^2 = \text{tr}((\widehat{W} - W)X X^\top (\widehat{W} - W)^\top) = \text{tr}((\widehat{W} - W)G_X(\widehat{W} - W)^\top).$$

□

**Corollary B.4** (Exact factorization of the conditional likelihood). *For  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L}),$  the conditional density of  $Y$  given  $X$  under the grouped law  $Q_\eta^{(n)}$  factorizes as*

$$p_\eta(Y | X) = g_X(Y) q_\eta(\widehat{W} | X), \quad (24)$$

where

$$g_X(Y) := (2\pi\sigma^2)^{-d(n-d)/2} |G_X|^{-d/2} \exp\left(-\frac{\|R\|_F^2}{2\sigma^2}\right) \quad (25)$$

is free of  $\eta$ , and

$$q_\eta(W | X) := (2\pi)^{-d^2/2} (\sigma^2)^{-d^2/2} |G_X|^{d/2} \sum_{z \in \mathcal{L}} \pi_z \exp\left(-\frac{1}{2\sigma^2} \text{tr}((W - A_z)G_X(W - A_z)^\top)\right) \quad (26)$$

is the conditional density of  $\widehat{W}$  given  $X$  under  $\eta$ .

In particular, conditional on  $X$ , the statistic  $\widehat{W}$  is sufficient for  $\eta$ .

*Proof.* Under  $\eta = (\pi, \Theta)$ ,

$$p_\eta(Y | X) = \sum_{z \in \mathcal{L}} \pi_z p_{A_z}(Y | X),$$

where

$$p_{A_z}(Y | X) = (2\pi\sigma^2)^{-dn/2} \exp\left(-\frac{\|Y - A_z X\|_F^2}{2\sigma^2}\right).$$

Applying Theorem B.3 with  $W = A_z$  yields

$$\|Y - A_z X\|_F^2 = \text{tr}((\widehat{W} - A_z)G_X(\widehat{W} - A_z)^\top) + \|R\|_F^2.$$

Substituting this into the display above and factoring out the term depending only on  $(X, Y)$  but not on  $\eta$  gives Equation (24) with  $g_X$  and  $q_\eta$  as in Equations (25) and (26). The sufficiency claim follows from the Neyman–Fisher factorization theorem applied conditionally on  $X$ .  $\square$

## B.2. Exact conditional law of the reduced statistic

We now compute the conditional law of  $\widehat{W}$  and the residual  $R$ , first under a fixed latent program and then under the full task-family law.

**Proposition B.5** (Exact matrix-normal reduction). *Fix  $z \in \mathcal{L}$ . Under the model*

$$Y = A_z X + E, \quad E = [\varepsilon_1, \dots, \varepsilon_n], \quad \varepsilon_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2 I_d),$$

the reduced statistics satisfy

$$\widehat{W} = A_z + EX^\top G_X^{-1}, \quad R = E(I_n - P_X). \quad (27)$$

Conditional on  $X$  and  $Z = z$ ,

$$\widehat{W} \sim \mathcal{MN}_{d,d}(A_z, \sigma^2 I_d, G_X^{-1}), \quad (28)$$

$$R \sim \mathcal{MN}_{d,n}(0, \sigma^2 I_d, I_n - P_X), \quad (29)$$

and  $\widehat{W}$  is independent of  $R$ .

*Proof.* Starting from  $Y = A_z X + E$ , the definition of  $\widehat{W}$  gives

$$\widehat{W} = (A_z X + E)X^\top G_X^{-1} = A_z X X^\top G_X^{-1} + EX^\top G_X^{-1} = A_z + EX^\top G_X^{-1},$$

which proves the first identity in Equation (27). Similarly,

$$R = (A_z X + E)(I_n - P_X) = A_z X(I_n - P_X) + E(I_n - P_X) = E(I_n - P_X),$$

because  $X(I_n - P_X) = 0$ . This proves the second identity in Equation (27).

Since  $\text{vec}(E) \sim \mathcal{N}(0, \sigma^2 I_{dn})$ , we have

$$\text{vec}(\widehat{W} - A_z) = \text{vec}(EX^\top G_X^{-1}) = ((G_X^{-1} X) \otimes I_d) \text{vec}(E).$$

Therefore,

$$\text{vec}(\widehat{W} - A_z) | X, Z = z \sim \mathcal{N}(0, \sigma^2 (G_X^{-1} X X^\top G_X^{-1}) \otimes I_d) = \mathcal{N}(0, \sigma^2 G_X^{-1} \otimes I_d),$$

which is exactly Equation (28).

Likewise,

$$\text{vec}(R) = \text{vec}(E(I_n - P_X)) = ((I_n - P_X) \otimes I_d) \text{vec}(E),$$

so

$$\text{vec}(R) \mid X, Z = z \sim \mathcal{N}(0, \sigma^2(I_n - P_X) \otimes I_d),$$

which gives Equation (29).

Finally, still conditional on  $X$  and  $Z = z$ , the pair  $(\widehat{W}, R)$  is jointly Gaussian because both are linear transformations of  $E$ . Its cross-covariance is

$$\begin{aligned} \text{Cov}(\text{vec}(\widehat{W} - A_z), \text{vec}(R) \mid X, Z = z) &= \sigma^2((G_X^{-1}X)(I_n - P_X)) \otimes I_d \\ &= 0, \end{aligned}$$

because  $X(I_n - P_X) = 0$ . Hence  $\widehat{W}$  and  $R$  are conditionally independent.  $\square$

**Corollary B.6** (Mixture representation under the task-family law). *Fix  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$ . Conditional on  $X$ , the reduced statistic  $\widehat{W}$  has the finite-mixture representation*

$$\mathcal{L}(\widehat{W} \mid X, \eta) = \sum_{z \in \mathcal{L}} \pi_z \mathcal{MN}_{d,d}(A_z, \sigma^2 I_d, G_X^{-1}). \quad (30)$$

Moreover, conditional on  $X$  and  $\eta$ , the residual  $R$  is independent of  $\widehat{W}$ , independent of the latent program  $Z$ , and has parameter-free law

$$R \mid X, \eta \sim \mathcal{MN}_{d,n}(0, \sigma^2 I_d, I_n - P_X).$$

*Proof.* The mixture representation Equation (30) follows immediately from Theorem B.5 by marginalizing over  $Z \sim \pi$ .

For the second claim, Theorem B.5 gives, for every fixed  $z \in \mathcal{L}$ ,

$$\mathcal{L}(R \mid X, Z = z, \eta) = \mathcal{MN}_{d,n}(0, \sigma^2 I_d, I_n - P_X),$$

which does not depend on  $z$  or on  $\eta$ . Hence  $R$  is conditionally independent of  $Z$  given  $X, \eta$ . Since  $R$  is also conditionally independent of  $\widehat{W}$  given  $X, Z, \eta$ , and its conditional law does not depend on  $Z$ , integrating out  $Z$  yields conditional independence of  $R$  and  $\widehat{W}$  given  $X, \eta$ .  $\square$

**Proposition B.7** (Exact posterior reduction). *Let  $\Pi$  be any prior on  $\Xi(B, \mathcal{L})$ . For each task  $i \in [m]$ , define*

$$\widehat{W}_i := Y_i X_i^\top (X_i X_i^\top)^{-1}.$$

Then for every measurable  $B \subseteq \Xi(B, \mathcal{L})$ ,

$$\Pi(B \mid X_{1:m}, Y_{1:m}) = \Pi(B \mid X_{1:m}, \widehat{W}_{1:m}) \quad \text{a.s.} \quad (31)$$

under the joint law induced by  $\Pi$  and the sampling model.

*Proof.* By independence across tasks and Theorem B.4, the conditional likelihood satisfies

$$p_\eta(Y_{1:m} \mid X_{1:m}) = \prod_{i=1}^m g_{X_i}(Y_i) \prod_{i=1}^m q_\eta(\widehat{W}_i \mid X_i).$$

The factor  $\prod_{i=1}^m g_{X_i}(Y_i)$  is free of  $\eta$ . Therefore Bayes' formula gives

$$\Pi(B \mid X_{1:m}, Y_{1:m}) = \frac{\int_B \prod_{i=1}^m q_\eta(\widehat{W}_i \mid X_i) \Pi(d\eta)}{\int_{\Xi(B, \mathcal{L})} \prod_{i=1}^m q_\eta(\widehat{W}_i \mid X_i) \Pi(d\eta)},$$

which is precisely the posterior based on  $(X_{1:m}, \widehat{W}_{1:m})$ . This proves Equation (31).  $\square$

**Corollary B.8** (Exact reduced experiment). *Under  $\eta = (\pi, \Theta)$ , conditional on the designs  $X_{1:m}$ , the reduced observations admit the representation*

$$Z_i \stackrel{\text{i.i.d.}}{\sim} \pi, \quad \widehat{W}_i = A_{Z_i} + \Xi_i, \quad \Xi_i | X_{1:m} \stackrel{\text{ind}}{\sim} \mathcal{MN}_{d,d}(0, \sigma^2 I_d, G_{X_i}^{-1}), \quad (32)$$

with  $\Xi_i$  conditionally independent of  $Z_i$  and independent across  $i \in [m]$ .

*Proof.* This is an immediate consequence of Theorem B.5 and independence across tasks.  $\square$

### B.3. Information geometry of the reduced experiment

The previous subsection shows that the grouped meta-learning problem reduces exactly to a finite mixture of heteroskedastic Gaussian kernels centered at the composed operators  $A_z$ . The next lemma records the resulting information geometry, which will be used later to construct tests and to relate structural perturbations in operator space to Hellinger and Kullback–Leibler separation.

For  $X$  with  $G_X$  invertible and  $W \in \mathbb{R}^{d \times d}$ , define

$$\mathcal{N}_X^W := \mathcal{MN}_{d,d}(W, \sigma^2 I_d, G_X^{-1}).$$

**Lemma B.9** (Exact conditional divergences). *For any full-row-rank design  $X$  and any  $W, W' \in \mathbb{R}^{d \times d}$ ,*

$$\begin{aligned} \text{KL}(\mathcal{N}_X^W \| \mathcal{N}_X^{W'}) &= \frac{1}{2\sigma^2} \text{tr}((W - W')G_X(W - W')^\top) \\ &= \frac{1}{2\sigma^2} \|(W - W')X\|_F^2, \end{aligned} \quad (33)$$

and

$$\begin{aligned} h^2(\mathcal{N}_X^W, \mathcal{N}_X^{W'}) &= 1 - \exp\left(-\frac{1}{8\sigma^2} \text{tr}((W - W')G_X(W - W')^\top)\right) \\ &= 1 - \exp\left(-\frac{1}{8\sigma^2} \|(W - W')X\|_F^2\right). \end{aligned} \quad (34)$$

*Proof.* Let  $\Delta := W - W'$ . By Theorem B.1,

$$\text{vec}(\mathcal{N}_X^W) \sim \mathcal{N}(\text{vec}(W), \Sigma_X), \quad \Sigma_X := G_X^{-1} \otimes \sigma^2 I_d.$$

The inverse covariance is

$$\Sigma_X^{-1} = G_X \otimes \sigma^{-2} I_d.$$

For Gaussian laws with common covariance, the exact formulas are

$$\text{KL}(\mathcal{N}(m, \Sigma) \| \mathcal{N}(m', \Sigma)) = \frac{1}{2} (m - m')^\top \Sigma^{-1} (m - m'),$$

and

$$h^2(\mathcal{N}(m, \Sigma), \mathcal{N}(m', \Sigma)) = 1 - \exp\left(-\frac{1}{8} (m - m')^\top \Sigma^{-1} (m - m')\right).$$

Applying these formulas with  $m = \text{vec}(W)$ ,  $m' = \text{vec}(W')$ , and  $\Sigma = \Sigma_X$ , we obtain

$$(\text{vec } \Delta)^\top \Sigma_X^{-1} \text{vec } \Delta = \frac{1}{\sigma^2} (\text{vec } \Delta)^\top (G_X \otimes I_d) \text{vec } \Delta = \frac{1}{\sigma^2} \text{tr}(\Delta G_X \Delta^\top).$$

Since  $G_X = XX^\top$ ,

$$\text{tr}(\Delta G_X \Delta^\top) = \text{tr}(\Delta X X^\top \Delta^\top) = \|\Delta X\|_F^2.$$

Substituting this identity into the Gaussian formulas yields Equations (33) and (34).  $\square$

#### B.4. Design concentration and effective noise scale

For later rate arguments, it is convenient to isolate a high-probability design event on which the reduced experiment behaves like a Gaussian location model with effective noise level  $n^{-1/2}$ .

For  $0 < \delta < 1$ , define the event

$$\mathcal{E}_n(\delta) := \{(1 - \delta)nI_d \preceq G_X \preceq (1 + \delta)nI_d\}. \quad (35)$$

**Lemma B.10** (Gaussian design concentration). *Fix  $0 < \delta < 1$ . There exist constants  $c_{d,\delta}, C_{d,\delta} \in (0, \infty)$ , depending only on  $d$  and  $\delta$ , such that*

$$\Pr(\mathcal{E}_n(\delta)^c) \leq C_{d,\delta} e^{-c_{d,\delta} n}. \quad (36)$$

The bound is uniform over  $\eta \in \Xi(B, \mathcal{L})$ , since the design law is independent of  $\eta$ .

*Proof.* Because the columns of  $X$  are i.i.d.  $\mathcal{N}(0, I_d)$ , the matrix  $n^{-1}G_X$  is the sample covariance matrix of  $n$  Gaussian observations in dimension  $d$ . The exponential tail bound Equation (36) is standard; see, for example, the discussion of Gaussian covariance concentration in Vershynin (2018, Chapter 4). Since the law of  $X$  does not depend on  $\eta$ , the bound is automatically uniform over the parameter space.  $\square$

**Corollary B.11** (Information-geometric sandwich on the good design event). *Fix  $0 < \delta < 1$ . On the event  $\mathcal{E}_n(\delta)$ , for every  $W, W' \in \mathbb{R}^{d \times d}$ ,*

$$\frac{(1 - \delta)n}{2\sigma^2} \|W - W'\|_F^2 \leq \text{KL}(\mathcal{N}_X^W \parallel \mathcal{N}_X^{W'}) \leq \frac{(1 + \delta)n}{2\sigma^2} \|W - W'\|_F^2, \quad (37)$$

and

$$\begin{aligned} 1 - \exp\left(-\frac{(1 - \delta)n}{8\sigma^2} \|W - W'\|_F^2\right) &\leq h^2(\mathcal{N}_X^W, \mathcal{N}_X^{W'}) \\ &\leq 1 - \exp\left(-\frac{(1 + \delta)n}{8\sigma^2} \|W - W'\|_F^2\right). \end{aligned} \quad (38)$$

In particular, whenever

$$\frac{(1 + \delta)n}{8\sigma^2} \|W - W'\|_F^2 \leq 1,$$

there exist constants  $0 < c_{\delta,\sigma} \leq C_{\delta,\sigma} < \infty$  such that

$$c_{\delta,\sigma} n \|W - W'\|_F^2 \leq h^2(\mathcal{N}_X^W, \mathcal{N}_X^{W'}) \leq C_{\delta,\sigma} n \|W - W'\|_F^2. \quad (39)$$

*Proof.* On  $\mathcal{E}_n(\delta)$ , the Gram matrix satisfies

$$(1 - \delta)nI_d \preceq G_X \preceq (1 + \delta)nI_d.$$

Hence for  $\Delta := W - W'$ ,

$$(1 - \delta)n \|\Delta\|_F^2 \leq \text{tr}(\Delta G_X \Delta^\top) \leq (1 + \delta)n \|\Delta\|_F^2.$$

Substituting into Equations (33) and (34) yields Equations (37) and (38).

For the local Hellinger bound, set

$$u := \frac{1}{8\sigma^2} \text{tr}(\Delta G_X \Delta^\top).$$

If  $u \leq 1$ , then the elementary inequalities

$$(1 - e^{-1})u \leq 1 - e^{-u} \leq u$$

imply

$$(1 - e^{-1})u \leq h^2(\mathcal{N}_X^W, \mathcal{N}_X^{W'}) \leq u.$$

Using the sandwich on  $u$  furnished by  $\mathcal{E}_n(\delta)$  gives Equation (39).  $\square$

*Remark B.12* (Role in later proofs). Theorems B.7 and B.8 show that every posterior statement in the paper may be proved in the reduced experiment  $(X_i, \widehat{W}_i)_{i=1}^m$ . Theorems B.9 and B.11 then convert operator perturbations into exact Hellinger and Kullback–Leibler separation. This is the technical bridge from latent task grammars to the testing and contraction arguments developed in the next appendices.

## C. Contextual Observability and Identifiability

Recent theory has made clear that compositional generalization is controlled by structural properties of the task family rather than by raw scale alone. In particular, support conditions matter for recovering reusable modules in teacher–student meta-learning (Schug et al., 2024), special autoregressive task families can permit generalization from  $\tilde{O}(D)$  tasks to a  $D^T$  family (Abedsoltan et al., 2025), and coverage / path-ambiguity considerations constrain what pattern-matching learners can infer from input–output pairs alone (Chang et al., 2025). The purpose of this appendix is to isolate the exact operator-level observability statements needed in our grouped Bayesian model.

The appendix has two logically distinct parts. First, we show that the pairwise quantity  $\Delta_{\text{ctx}}$  introduced in the main text is an exact information measure for one-hole substitutions and yields a sharp impossibility theorem in the single-occurrence regime. Second, we show that  $\Delta_{\text{ctx}}$  by itself is not strong enough for full local identification of the latent grammar. For that stronger conclusion we introduce a *linearized contextual observability* modulus  $\lambda_{\text{lin}}$ , prove local injectivity of the program-operator map, and then derive local quotient identifiability of the full parameter.

### C.1. One-hole substitutions and exact information separation

We begin by formalizing the operator-level effect of substituting one primitive module for another inside a fixed context.

For strings  $u \in [r]^s$ ,  $v \in [r]^q$ , and  $a \in [r]$ , recall that

$$A_{uav} = A_v A_a A_u, \quad A_\emptyset = I_d.$$

For any full-row-rank design  $X \in \mathbb{R}^{d \times n}$ , define

$$\mathcal{N}_X^W := \mathcal{MN}_{d,d}(W, \sigma^2 I_d, G_X^{-1}), \quad G_X := X X^\top.$$

**Lemma C.1** (Exact one-hole substitution divergences). *Let  $u \in [r]^s$ ,  $v \in [r]^q$ , and  $a, b \in [r]$  be distinct. For every full-row-rank design  $X$  and every  $n \geq d + 1$ ,*

$$\text{KL}\left(\mathcal{N}_X^{A_{uav}} \parallel \mathcal{N}_X^{A_{ubv}}\right) = \frac{1}{2\sigma^2} \|A_v(A_a - A_b)A_u X\|_F^2, \quad (40)$$

$$h^2\left(\mathcal{N}_X^{A_{uav}}, \mathcal{N}_X^{A_{ubv}}\right) = 1 - \exp\left(-\frac{1}{8\sigma^2} \|A_v(A_a - A_b)A_u X\|_F^2\right). \quad (41)$$

*Proof.* Apply Theorem B.9 from Appendix B with

$$W = A_{uav} = A_v A_a A_u, \quad W' = A_{ubv} = A_v A_b A_u.$$

Then  $W - W' = A_v(A_a - A_b)A_u$ , and substituting into Equations (33) and (34) yields Equations (40) and (41).  $\square$

The next proposition shows that  $\Delta_{\text{ctx}}$  is not merely a heuristic notion of difficulty: it is exactly the  $\nu_\pi$ -average information for distinguishing one primitive module from another through random contexts induced by the task-family law.

**Proposition C.2** (Context-averaged information equals contextual observability). *Fix  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$  and distinct  $a, b \in [r]$ . Let  $X$  be a full-row-rank design. Then*

$$\mathbb{E}_{(u,v) \sim \nu_\pi} \text{KL}\left(\mathcal{N}_X^{A_{uav}} \parallel \mathcal{N}_X^{A_{ubv}}\right) = \frac{1}{2\sigma^2} \sum_{(u,v) \in \mathcal{C}_\mathcal{L}} \nu_\pi(u, v) \|A_v(A_a - A_b)A_u X\|_F^2. \quad (42)$$

Moreover, on the good design event  $\mathcal{E}_n(\delta)$  of Equation (35),

$$\frac{(1 - \delta)n}{2\sigma^2} \Delta_{\text{ctx}}(a, b; \eta) \leq \mathbb{E}_{(u,v) \sim \nu_\pi} \text{KL}\left(\mathcal{N}_X^{A_{uav}} \parallel \mathcal{N}_X^{A_{ubv}}\right) \leq \frac{(1 + \delta)n}{2\sigma^2} \Delta_{\text{ctx}}(a, b; \eta). \quad (43)$$

Finally, if in addition

$$\frac{(1 + \delta)n}{8\sigma^2} \max_{(u,v) \in \mathcal{C}_\mathcal{L}} \|A_v(A_a - A_b)A_u\|_F^2 \leq 1,$$

then there exist constants  $0 < c_{\delta, \sigma} \leq C_{\delta, \sigma} < \infty$  such that

$$c_{\delta, \sigma} n \Delta_{\text{ctx}}(a, b; \eta) \leq \mathbb{E}_{(u,v) \sim \nu_\pi} h^2\left(\mathcal{N}_X^{A_{uav}}, \mathcal{N}_X^{A_{ubv}}\right) \leq C_{\delta, \sigma} n \Delta_{\text{ctx}}(a, b; \eta). \quad (44)$$

*Proof.* The exact identity Equation (42) follows by multiplying Equation (40) by  $\nu_\pi(u, v)$  and summing over  $(u, v) \in \mathcal{C}_\mathcal{L}$ .

On  $\mathcal{E}_n(\delta)$ , the Gram matrix satisfies

$$(1 - \delta)nI_d \preceq G_X \preceq (1 + \delta)nI_d.$$

Hence for every  $(u, v) \in \mathcal{C}_\mathcal{L}$ ,

$$(1 - \delta)n\|A_v(A_a - A_b)A_u\|_F^2 \leq \|A_v(A_a - A_b)A_uX\|_F^2 \leq (1 + \delta)n\|A_v(A_a - A_b)A_u\|_F^2.$$

Substituting this bound into Equation (42) yields Equation (43).

For the Hellinger claim, apply Equation (41) together with the local inequality

$$(1 - e^{-1})u \leq 1 - e^{-u} \leq u \quad \text{for } u \in [0, 1].$$

Under the displayed small-separation condition, the exponent in Equation (41) is uniformly at most 1 on  $\mathcal{E}_n(\delta)$ . Averaging over  $(u, v) \sim \nu_\pi$  then gives Equation (44).  $\square$

*Remark C.3* (What  $\Delta_{\text{ctx}}$  does and does not control). Theorem C.2 shows that  $\Delta_{\text{ctx}}$  is the exact pairwise information for single-site substitutions through random task-family contexts. This is the operator-level analogue of the data-centric support / coverage ideas emphasized in recent compositional-generalization theory (Chang et al., 2025; Schug et al., 2024). However,  $\Delta_{\text{ctx}}$  is a *pairwise* quantity. It controls first-order distinguishability of one primitive from another, but it does not by itself guarantee full injectivity of the entire map  $\Theta \mapsto (A_z)_{z \in \mathcal{L}}$ . That stronger conclusion requires the linearized modulus introduced in Theorem C.9 below.

## C.2. Zero contextual observability and exact collapse in the single-occurrence regime

The impossibility statement for  $\Delta_{\text{ctx}} = 0$  is exact in a natural and important special case: languages in which no primitive symbol appears more than once in a program. This covers many synthetic task families used to study compositional generalization, including settings where a task is an ordered composition of distinct modules.

**Definition C.4** (Single-occurrence language). We say that the admissible language  $\mathcal{L}$  is *single-occurrence* if for every  $z = (z_1, \dots, z_t) \in \mathcal{L}$  and every  $a \in [r]$ ,

$$\#\{\ell \in \{1, \dots, t\} : z_\ell = a\} \leq 1.$$

**Proposition C.5** (Exact collapse at zero contextual observability in single-occurrence languages). *Assume that  $\mathcal{L}$  is single-occurrence, and fix  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$ . Let  $a, b \in [r]$  be distinct and suppose that*

$$\Delta_{\text{ctx}}(a, b; \eta) = 0.$$

Define the modified library

$$\Theta^{(a \leftarrow b)} := (A_1, \dots, A_{b-1}, A_a, A_{b+1}, \dots, A_r),$$

and let

$$\tilde{\eta} := (\pi, \Theta^{(a \leftarrow b)}).$$

Then

$$A_z(\Theta^{(a \leftarrow b)}) = A_z(\Theta) \quad \text{for every } z \in \mathcal{L}, \quad (45)$$

and therefore

$$Q_{\tilde{\eta}}^{(n)} = Q_{\eta}^{(n)} \quad \text{for every } n \geq 1. \quad (46)$$

*Proof.* Since  $\Delta_{\text{ctx}}(a, b; \eta) = 0$  and every term in its definition is nonnegative, we must have

$$A_v(A_a - A_b)A_u = 0 \quad \text{for every } (u, v) \in \mathcal{C}_\mathcal{L} \text{ with } \nu_\pi(u, v) > 0. \quad (47)$$

Because  $\pi \in \mathfrak{P}_\mathcal{L}^\circ$ , every program  $z \in \mathcal{L}$  has strictly positive probability  $\pi_z > 0$ . Hence every actual occurrence of a symbol in every  $z \in \mathcal{L}$  contributes strictly positive mass to the context law  $\nu_\pi$ .

Fix  $z \in \mathcal{L}$ . If  $z$  contains no occurrence of  $b$ , then  $A_z(\Theta^{(a \leftarrow b)}) = A_z(\Theta)$  trivially. If  $z$  contains one occurrence of  $b$ , then by the single-occurrence assumption it can be written uniquely as

$$z = ubv$$

for some  $u, v$  with  $(u, v) = \text{ctx}_\ell(z)$  at the unique position  $\ell$  carrying  $b$ . Since  $\pi_z > 0$ , the context  $(u, v)$  has  $\nu_\pi(u, v) > 0$ , so Equation (47) yields

$$A_v(A_a - A_b)A_u = 0.$$

Therefore

$$A_z(\Theta^{(a \leftarrow b)}) - A_z(\Theta) = A_v A_a A_u - A_v A_b A_u = A_v(A_a - A_b)A_u = 0.$$

This proves Equation (45). Summing the identical component laws with the same weights  $\pi_z$  gives Equation (46).  $\square$

*Remark C.6* (Why the single-occurrence restriction is real). Without the single-occurrence restriction,  $\Delta_{\text{ctx}}(a, b; \eta) = 0$  only rules out *first-order* contextual visibility of the substitution  $b \mapsto a$ . If a program contains multiple occurrences of  $b$ , then higher-order replacement terms can survive even when every one-hole context is individually unobservable. In repeated-symbol languages, the correct exact obstruction is therefore not the scalar  $\Delta_{\text{ctx}}$  alone, but the full linearized operator introduced next.

### C.3. A stronger local observability modulus

We now introduce the quantity that governs full local identifiability of the module library from the induced family of composed operators.

**Definition C.7** (Program-operator map and library norm). Let

$$\Phi : \mathbb{A}_B^r \rightarrow (\mathbb{R}^{d \times d})^{\mathcal{L}}, \quad \Phi(\Theta) := (A_z(\Theta))_{z \in \mathcal{L}}.$$

For  $H = (H_1, \dots, H_r) \in (\mathbb{R}^{d \times d})^r$ , define

$$\|H\|_{\text{lib}}^2 := \sum_{a=1}^r \|H_a\|_F^2.$$

For  $U = (U_z)_{z \in \mathcal{L}} \in (\mathbb{R}^{d \times d})^{\mathcal{L}}$  and  $\pi \in \mathfrak{P}_{\mathcal{L}}^\circ$ , define

$$\|U\|_\pi^2 := \sum_{z \in \mathcal{L}} \pi_z \|U_z\|_F^2.$$

For  $z = (z_1, \dots, z_t) \in \mathcal{L}$  and  $\ell \in \{1, \dots, t\}$ , write

$$z_{<\ell} := (z_1, \dots, z_{\ell-1}), \quad z_{>\ell} := (z_{\ell+1}, \dots, z_t).$$

**Lemma C.8** (Derivative of the program-operator map). *The map  $\Phi$  is polynomial, hence  $C^\infty$ , and its Fréchet derivative at  $\Theta = (A_1, \dots, A_r)$  is given by*

$$[D\Phi_\Theta(H)]_z = \sum_{\ell=1}^{|z|} A_{z_{>\ell}} H_{z_\ell} A_{z_{<\ell}} \quad \text{for every } z \in \mathcal{L}. \quad (48)$$

*Proof.* Fix  $z = (z_1, \dots, z_t)$ . The map

$$\Theta \mapsto A_z(\Theta) = A_{z_t} \cdots A_{z_1}$$

is a noncommutative polynomial of degree  $t \leq L$  in the matrix coordinates  $(A_1, \dots, A_r)$ , hence smooth. Differentiating the product in each location where a perturbation may enter yields

$$D(A_z)_\Theta(H) = \sum_{\ell=1}^t A_{z_t} \cdots A_{z_{\ell+1}} H_{z_\ell} A_{z_{\ell-1}} \cdots A_{z_1},$$

which is exactly Equation (48). Since  $\mathcal{L}$  is finite, the statement follows coordinatewise.  $\square$

**Definition C.9** (Linearized contextual observability). For  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$ , define the *linearized contextual observability modulus*

$$\lambda_{\text{lin}}(\eta) := \inf_{\|H\|_{\text{lib}}=1} \|D\Phi_{\Theta}(H)\|_{\pi}^2. \quad (49)$$

Equivalently,  $\lambda_{\text{lin}}(\eta)$  is the smallest singular value squared of the linear map

$$D\Phi_{\Theta} : ((\mathbb{R}^{d \times d})^r, \|\cdot\|_{\text{lib}}) \rightarrow ((\mathbb{R}^{d \times d})^{\mathcal{L}}, \|\cdot\|_{\pi}).$$

*Remark C.10* (Why  $\lambda_{\text{lin}}$  is stronger than  $\Delta_{\text{ctx}}$ ). The scalar  $\Delta_{\text{ctx}}$  probes pairwise visibility of replacing one primitive by another inside a random one-hole context. By contrast,  $\lambda_{\text{lin}}$  controls *all* infinitesimal perturbations of the module library simultaneously. This distinction is unavoidable: full local identification of the map  $\Theta \mapsto (A_z)_{z \in \mathcal{L}}$  requires ruling out arbitrary coordinated cancellations, not only pairwise substitutions.

The next lemma supplies the quadratic remainder estimate needed to convert injectivity of the derivative into local injectivity of the nonlinear program-operator map.

**Lemma C.11** (Quadratic Taylor remainder). Fix  $\eta^* = (\pi^*, \Theta^*) \in \Xi(B, \mathcal{L})$ . There exists a constant

$$C_{\text{quad}} = C_{\text{quad}}(\mathcal{L}, r, d, L, B, \pi^*) < \infty$$

such that for every  $H \in (\mathbb{R}^{d \times d})^r$  with  $\Theta^* + H \in \mathbb{A}_B^r$ ,

$$\|\Phi(\Theta^* + H) - \Phi(\Theta^*) - D\Phi_{\Theta^*}(H)\|_{\pi^*} \leq C_{\text{quad}} \|H\|_{\text{lib}}^2. \quad (50)$$

*Proof.* Each coordinate map  $\Theta \mapsto A_z(\Theta)$  is a polynomial of degree at most  $L$  on the finite-dimensional Euclidean space  $(\mathbb{R}^{d \times d})^r$ . Therefore its second derivative is continuous. Since  $\mathbb{A}_B^r$  is compact, the second derivative of each coordinate is uniformly bounded on  $\mathbb{A}_B^r$ . Summing over the finitely many coordinates  $z \in \mathcal{L}$  and using the weighted norm  $\|\cdot\|_{\pi^*}$  shows that the Hessian of  $\Phi$  is uniformly bounded on  $\mathbb{A}_B^r$  as a bilinear map from  $(\mathbb{R}^{d \times d})^r \times (\mathbb{R}^{d \times d})^r$  into  $((\mathbb{R}^{d \times d})^{\mathcal{L}}, \|\cdot\|_{\pi^*})$ . The finite-dimensional Taylor theorem with integral remainder then yields Equation (50).  $\square$

**Proposition C.12** (Local injectivity of the program-operator map). Fix  $\eta^* = (\pi^*, \Theta^*) \in \Xi(B, \mathcal{L})$  and suppose that

$$\lambda_{\text{lin}}(\eta^*) > 0.$$

Then there exist constants  $\rho_{\text{lin}} > 0$  and  $c_{\text{lin}} > 0$  such that for every  $\Theta \in \mathbb{A}_B^r$  with

$$\|\Theta - \Theta^*\|_{\text{lib}} \leq \rho_{\text{lin}},$$

one has

$$\|\Phi(\Theta) - \Phi(\Theta^*)\|_{\pi^*} \geq c_{\text{lin}} \|\Theta - \Theta^*\|_{\text{lib}}. \quad (51)$$

In particular,  $\Phi$  is locally injective at  $\Theta^*$ .

*Proof.* Set

$$\lambda_{\star} := \lambda_{\text{lin}}(\eta^*) > 0, \quad c_{\star} := \sqrt{\lambda_{\star}}.$$

By definition of  $\lambda_{\text{lin}}$ ,

$$\|D\Phi_{\Theta^*}(H)\|_{\pi^*} \geq c_{\star} \|H\|_{\text{lib}} \quad \text{for every } H \in (\mathbb{R}^{d \times d})^r.$$

Now let  $H := \Theta - \Theta^*$ . By Theorem C.11,

$$\|\Phi(\Theta) - \Phi(\Theta^*)\|_{\pi^*} \geq \|D\Phi_{\Theta^*}(H)\|_{\pi^*} - C_{\text{quad}} \|H\|_{\text{lib}}^2 \geq c_{\star} \|H\|_{\text{lib}} - C_{\text{quad}} \|H\|_{\text{lib}}^2.$$

Choose

$$\rho_{\text{lin}} := \min \left\{ 1, \frac{c_{\star}}{2C_{\text{quad}}} \right\}, \quad c_{\text{lin}} := \frac{c_{\star}}{2}.$$

Then for  $\|H\|_{\text{lib}} \leq \rho_{\text{lin}}$ ,

$$C_{\text{quad}} \|H\|_{\text{lib}}^2 \leq \frac{c_{\star}}{2} \|H\|_{\text{lib}},$$

whence

$$\|\Phi(\Theta) - \Phi(\Theta^*)\|_{\pi^*} \geq \frac{c_{\star}}{2} \|H\|_{\text{lib}} = c_{\text{lin}} \|\Theta - \Theta^*\|_{\text{lib}}.$$

This is Equation (51).  $\square$

#### C.4. From grouped laws to the induced mixing measure

The grouped law  $Q_\eta^{(n)}$  depends on the latent grammar only through the induced discrete measure on composed operators.

**Definition C.13** (Induced mixing measure and component separation). For  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$ , define the induced finite measure on operator space

$$\mu_\eta := \sum_{z \in \mathcal{L}} \pi_z \delta_{A_z(\Theta)} \quad \text{on } \mathbb{R}^{d \times d}. \quad (52)$$

For  $\eta^* = (\pi^*, \Theta^*)$ , define the component-separation margin

$$\gamma_{\text{comp}}(\eta^*) := \min_{z \neq z'} \|A_z^* - A_{z'}^*\|_F. \quad (53)$$

*Remark C.14* (Mixture identifiability layer). Classical results of Teicher (1963); Yakowitz & Spragins (1968) show that finite mixtures of Gaussian kernels are identifiable. In our setting we can prove the relevant statement directly, because the reduced experiment of Appendix B is a finite mixture of matrix-normal kernels with a common covariance depending only on the design.

**Lemma C.15** (Injectivity of the common-covariance matrix-normal kernel). Fix a full-row-rank design  $X \in \mathbb{R}^{d \times n}$ . Let  $\mu$  and  $\nu$  be finite signed Borel measures on  $\mathbb{R}^{d \times d}$ . Suppose that the corresponding mixture densities satisfy

$$\int \varphi_X(W_{\text{obs}} - W) \mu(dW) = \int \varphi_X(W_{\text{obs}} - W) \nu(dW) \quad \text{for Lebesgue-a.e. } W_{\text{obs}} \in \mathbb{R}^{d \times d},$$

where  $\varphi_X(\cdot - W)$  denotes the density of  $\mathcal{MN}_{d,d}(W, \sigma^2 I_d, G_X^{-1})$ . Then  $\mu = \nu$ .

*Proof.* Vectorize the matrices. Let

$$p := d^2, \quad \Sigma_X := G_X^{-1} \otimes \sigma^2 I_d \in \mathbb{R}^{p \times p},$$

which is positive definite because  $G_X$  is. Let  $\bar{\mu}$  and  $\bar{\nu}$  be the pushforwards of  $\mu$  and  $\nu$  under the map  $W \mapsto \text{vec}(W)$ . Then the assumed equality becomes

$$\phi_{\Sigma_X} * \bar{\mu} = \phi_{\Sigma_X} * \bar{\nu} \quad \text{a.e. on } \mathbb{R}^p,$$

where  $\phi_{\Sigma_X}$  is the  $p$ -variate Gaussian density with covariance  $\Sigma_X$ , and  $*$  denotes convolution.

Taking Fourier transforms gives

$$\widehat{\phi_{\Sigma_X}}(t) \widehat{\bar{\mu}}(t) = \widehat{\phi_{\Sigma_X}}(t) \widehat{\bar{\nu}}(t) \quad \text{for all } t \in \mathbb{R}^p.$$

Since

$$\widehat{\phi_{\Sigma_X}}(t) = \exp\left(-\frac{1}{2} t^\top \Sigma_X t\right)$$

is strictly positive for every  $t$ , it follows that

$$\widehat{\bar{\mu}}(t) = \widehat{\bar{\nu}}(t) \quad \text{for all } t \in \mathbb{R}^p.$$

By the uniqueness theorem for finite Borel measures on  $\mathbb{R}^p$ ,  $\bar{\mu} = \bar{\nu}$ . Pushing back through the vectorization map yields  $\mu = \nu$ .  $\square$

**Proposition C.16** (The grouped law determines the induced mixing measure). Let  $\eta, \eta' \in \Xi(B, \mathcal{L})$ , and let  $n \geq d + 1$ . If

$$Q_\eta^{(n)} = Q_{\eta'}^{(n)}$$

as probability laws on one task dataset  $D = (X, Y)$ , then

$$\mu_\eta = \mu_{\eta'}$$

as finite measures on  $\mathbb{R}^{d \times d}$ .

*Proof.* Since  $\widehat{W} = YX^\top(XX^\top)^{-1}$  is a measurable function of  $D = (X, Y)$  on the almost-sure event  $\{XX^\top \text{ invertible}\}$ , equality of the grouped laws implies equality of the joint laws of  $(X, \widehat{W})$  under  $\eta$  and  $\eta'$ . By disintegration, for almost every  $X$  with full row rank,

$$\mathcal{L}(\widehat{W} \mid X, \eta) = \mathcal{L}(\widehat{W} \mid X, \eta').$$

By Theorem B.6 from Appendix B,

$$\mathcal{L}(\widehat{W} \mid X, \eta) = \int \mathcal{MN}_{d,d}(W, \sigma^2 I_d, G_X^{-1}) \mu_\eta(dW),$$

and analogously for  $\eta'$ . Applying Theorem C.15 for any such full-row-rank  $X$  yields  $\mu_\eta = \mu_{\eta'}$ .  $\square$

### C.5. Local quotient identifiability

We can now state and prove the exact local identifiability theorem.

**Theorem C.17** (Local quotient identifiability). *Fix  $\eta^* = (\pi^*, \Theta^*) \in \Xi(B, \mathcal{L})$ . Assume that*

$$\lambda_{\text{lin}}(\eta^*) > 0, \quad \gamma_{\text{comp}}(\eta^*) > 0. \quad (54)$$

*Then there exists  $\rho_{\text{id}} > 0$  such that the following holds. For any  $n \geq d + 1$  and any  $\eta \in \Xi(B, \mathcal{L})$ ,*

$$d_q(\eta, \eta^*) \leq \rho_{\text{id}} \quad \text{and} \quad Q_\eta^{(n)} = Q_{\eta^*}^{(n)}$$

*imply*

$$d_q(\eta, \eta^*) = 0.$$

*Equivalently, in a sufficiently small quotient neighborhood of  $\eta^*$ , the grouped law uniquely identifies the latent task-family parameter.*

*Proof.* Let

$$\gamma_* := \gamma_{\text{comp}}(\eta^*) > 0.$$

By continuity of  $\Phi$ , there exists  $\rho_{\text{sep}} > 0$  such that whenever

$$\|\Theta - \Theta^*\|_{\text{lib}} \leq \rho_{\text{sep}},$$

one has

$$\max_{z \in \mathcal{L}} \|A_z(\Theta) - A_z^*\|_F < \frac{\gamma_*}{4}. \quad (55)$$

Let  $\rho_{\text{lin}} > 0$  and  $c_{\text{lin}} > 0$  be given by Theorem C.12. Set

$$\rho_{\text{id}} := \min\{\rho_{\text{sep}}, \rho_{\text{lin}}\}.$$

Now fix  $\eta = (\pi, \Theta)$  with  $d_q(\eta, \eta^*) \leq \rho_{\text{id}}$ , and assume

$$Q_\eta^{(n)} = Q_{\eta^*}^{(n)}.$$

By definition of  $d_q$ , there exists  $\sigma \in \text{Aut}(\mathcal{L})$  such that

$$d_{\text{par}}(\sigma \cdot \eta, \eta^*) = d_q(\eta, \eta^*) \leq \rho_{\text{id}}.$$

Since grouped laws are quotient-invariant by Theorem A.4,

$$Q_{\sigma \cdot \eta}^{(n)} = Q_\eta^{(n)} = Q_{\eta^*}^{(n)}.$$

Hence it suffices to prove  $\sigma \cdot \eta = \eta^*$ . Replacing  $\eta$  by  $\sigma \cdot \eta$ , we may therefore assume from the outset that

$$d_{\text{par}}(\eta, \eta^*) \leq \rho_{\text{id}}.$$

By Theorem C.16,

$$\mu_\eta = \mu_{\eta^*}.$$

Write

$$\mu_{\eta^*} = \sum_{z \in \mathcal{L}} \pi_z^* \delta_{A_z^*}, \quad \mu_\eta = \sum_{z \in \mathcal{L}} \pi_z \delta_{A_z(\Theta)}.$$

By Equation (55), every atom  $A_z(\Theta)$  lies inside the open ball

$$B_z := B\left(A_z^*, \frac{\gamma_\star}{4}\right).$$

These balls are pairwise disjoint because the centers are separated by at least  $\gamma_\star$ . Moreover, if  $z' \neq z$ , then

$$\|A_{z'}(\Theta) - A_z^*\|_F \geq \|A_{z'}^* - A_z^*\|_F - \|A_{z'}(\Theta) - A_{z'}^*\|_F > \gamma_\star - \frac{\gamma_\star}{4} = \frac{3\gamma_\star}{4},$$

so  $A_{z'}(\Theta) \notin B_z$ . Hence  $B_z$  contains exactly one atom of  $\mu_\eta$ , namely  $A_z(\Theta)$ , and exactly one atom of  $\mu_{\eta^*}$ , namely  $A_z^*$ .

Since  $\mu_\eta = \mu_{\eta^*}$ , the two measures assign the same mass to  $B_z$ . Therefore

$$\pi_z = \mu_\eta(B_z) = \mu_{\eta^*}(B_z) = \pi_z^* \quad \text{for every } z \in \mathcal{L}.$$

Moreover, because the only atom of  $\mu_\eta$  inside  $B_z$  is  $A_z(\Theta)$ , and the only atom of  $\mu_{\eta^*}$  inside  $B_z$  is  $A_z^*$ , equality of the two atomic measures forces

$$A_z(\Theta) = A_z^* \quad \text{for every } z \in \mathcal{L}.$$

Thus

$$\Phi(\Theta) = \Phi(\Theta^*).$$

Since  $\|\Theta - \Theta^*\|_{\text{lib}} \leq \rho_{\text{id}} \leq \rho_{\text{lin}}$ , Theorem C.12 implies

$$0 = \|\Phi(\Theta) - \Phi(\Theta^*)\|_{\pi^*} \geq c_{\text{lin}} \|\Theta - \Theta^*\|_{\text{lib}},$$

hence  $\Theta = \Theta^*$ . Together with  $\pi = \pi^*$ , this shows  $\eta = \eta^*$ . Undoing the initial quotient alignment yields  $d_q(\eta, \eta^*) = 0$ .  $\square$

*Remark C.18 (Interpretation of the theorem).* The theorem cleanly separates two layers of identifiability. The first layer is *mixture identifiability*: the grouped law determines the induced discrete measure  $\mu_\eta$  on composed operators. The second is *grammar identifiability*: the operator family  $(A_z)_{z \in \mathcal{L}}$  locally determines the module library  $\Theta$ . The first layer is classical in spirit and follows from the common-covariance Gaussian kernel; the second is the genuinely compositional part, and it is governed by  $\lambda_{\text{lin}}$ , not by  $\Delta_{\text{ctx}}$  alone.

## D. Testing, Entropy, and Prior Mass

This appendix verifies the three technical ingredients needed for posterior contraction in the reduced experiment of Appendix B: exponentially consistent tests, metric entropy bounds, and local prior thickness. The overall architecture is standard in Bayesian nonparametrics (Ghosal et al., 2000), but the substantive content here is model-specific: we compute all bounds in the grouped stochastic-grammar experiment and make explicit how the within-task sample size  $n$  sharpens the geometry of the module parameters. For robust tests between Hellinger balls, we appeal to the Le Cam–Birgé theory in the form summarized by Birgé (2013).

### D.1. Conditional product experiment and predictive radius

Fix  $m \geq 1$  and deterministic designs  $X_{1:m} = (X_1, \dots, X_m)$  with each  $X_i \in \mathbb{R}^{d \times n}$  full row rank. For  $\eta \in \Xi(B, \mathcal{L})$ , let

$$q_{\eta,i}(\cdot) := q_\eta(\cdot \mid X_i)$$

denote the conditional density of the reduced statistic  $\widehat{W}_i$  given  $X_i$ , as defined in Equation (26). The conditional product law of the reduced experiment is

$$\mathbb{P}_{\eta,X}^{(m)} := \bigotimes_{i=1}^m q_{\eta,i}.$$

We write

$$\rho(P, Q) := \int \sqrt{pq} d\mu = 1 - h^2(P, Q)$$

for the Hellinger affinity of two dominated probability laws  $P$  and  $Q$ .

**Definition D.1** (Average conditional Hellinger metric). For  $\eta, \eta' \in \Xi(B, \mathcal{L})$ , define

$$\bar{h}_X^2(\eta, \eta') := \frac{1}{m} \sum_{i=1}^m h^2(q_{\eta, i}, q_{\eta', i}). \quad (56)$$

The predictive contraction scale of the grouped law is

$$\delta_{m,n}^2 := \frac{K_{\text{eff}} \log m + rd^2 \log(mn)}{m}. \quad (57)$$

This is the scale at which the entropy and prior-thickness bounds balance. The sharper  $1/(mn)$  scaling for the shared module library will emerge only after converting grouped-law contraction into structural error in Appendix F.

## D.2. Conditional Hellinger tests

We begin with the simple-vs-simple testing bound. It is included explicitly because it is the quantitative hinge on which the later net argument turns.

**Lemma D.2** (Likelihood-ratio test for the conditional product experiment). Fix  $\eta_0, \eta_1 \in \Xi(B, \mathcal{L})$  and deterministic designs  $X_{1:m}$ . There exists a measurable test  $\phi = \phi_{\eta_0, \eta_1, X} \in [0, 1]$  such that

$$\mathbb{E}_{\eta_0}[\phi \mid X_{1:m}] + \mathbb{E}_{\eta_1}[1 - \phi \mid X_{1:m}] \leq \rho\left(\mathbb{P}_{\eta_0, X}^{(m)}, \mathbb{P}_{\eta_1, X}^{(m)}\right) \leq \exp(-m \bar{h}_X^2(\eta_0, \eta_1)). \quad (58)$$

*Proof.* Let  $p_0$  and  $p_1$  denote the densities of  $\mathbb{P}_{\eta_0, X}^{(m)}$  and  $\mathbb{P}_{\eta_1, X}^{(m)}$  with respect to a common dominating measure on  $(\mathbb{R}^{d \times d})^m$ . Consider the likelihood-ratio test

$$\phi := \mathbf{1}\{p_1 \geq p_0\}.$$

Then

$$\mathbb{E}_{\eta_0}[\phi \mid X_{1:m}] + \mathbb{E}_{\eta_1}[1 - \phi \mid X_{1:m}] = \int \min\{p_0, p_1\} \leq \int \sqrt{p_0 p_1} = \rho\left(\mathbb{P}_{\eta_0, X}^{(m)}, \mathbb{P}_{\eta_1, X}^{(m)}\right).$$

Because the product experiment is conditionally independent across tasks,

$$\rho\left(\mathbb{P}_{\eta_0, X}^{(m)}, \mathbb{P}_{\eta_1, X}^{(m)}\right) = \prod_{i=1}^m \rho(q_{\eta_0, i}, q_{\eta_1, i}) = \prod_{i=1}^m (1 - h^2(q_{\eta_0, i}, q_{\eta_1, i})).$$

Using  $1 - u \leq e^{-u}$  for  $u \in [0, 1]$  gives

$$\prod_{i=1}^m (1 - h^2(q_{\eta_0, i}, q_{\eta_1, i})) \leq \exp\left(-\sum_{i=1}^m h^2(q_{\eta_0, i}, q_{\eta_1, i})\right) = \exp(-m \bar{h}_X^2(\eta_0, \eta_1)).$$

□

The next proposition upgrades Theorem D.2 to composite alternatives by combining a finite net with the classical Le Cam–Birgé robust test between Hellinger balls.

**Proposition D.3** (Composite tests from Hellinger coverings). Fix  $\eta^* \in \Xi(B, \mathcal{L})$ , deterministic designs  $X_{1:m}$ , and  $\varepsilon > 0$ . Let  $A \subseteq \Xi(B, \mathcal{L})$  satisfy

$$\inf_{\eta \in A} \bar{h}_X(\eta, \eta^*) \geq 4\varepsilon.$$

Then there exists a measurable test  $\phi_{A, \varepsilon, X} \in [0, 1]$  such that

$$\mathbb{E}_{\eta^*}[\phi_{A, \varepsilon, X} \mid X_{1:m}] \leq N(\varepsilon, A, \bar{h}_X) e^{-m\varepsilon^2}, \quad (59)$$

$$\sup_{\eta \in A} \mathbb{E}_{\eta}[1 - \phi_{A, \varepsilon, X} \mid X_{1:m}] \leq e^{-m\varepsilon^2}, \quad (60)$$

where  $N(\varepsilon, A, \bar{h}_X)$  denotes the  $\varepsilon$ -covering number of  $A$  under  $\bar{h}_X$ .

*Proof.* Let  $\{\eta_1, \dots, \eta_N\} \subseteq A$  be an  $\varepsilon$ -net of  $A$  with  $N = N(\varepsilon, A, \bar{h}_X)$ . Since every center  $\eta_j$  satisfies  $\bar{h}_X(\eta_j, \eta^*) \geq 4\varepsilon$ , the robust Hellinger-ball testing theorem of Le Cam–Birgé yields, for each  $j$ , a test  $\phi_j$  such that

$$\sup_{\bar{h}_X(\eta, \eta^*) \leq \varepsilon} \mathbb{E}_\eta[\phi_j | X_{1:m}] \leq e^{-m\varepsilon^2}, \quad \sup_{\bar{h}_X(\eta, \eta_j) \leq \varepsilon} \mathbb{E}_\eta[1 - \phi_j | X_{1:m}] \leq e^{-m\varepsilon^2}.$$

See Birgé (2013) for a modern statement encompassing independent non-identically distributed experiments.

Define

$$\phi_{A, \varepsilon, X} := \max_{1 \leq j \leq N} \phi_j.$$

Then

$$\mathbb{E}_{\eta^*}[\phi_{A, \varepsilon, X} | X_{1:m}] \leq \sum_{j=1}^N \mathbb{E}_{\eta^*}[\phi_j | X_{1:m}] \leq Ne^{-m\varepsilon^2},$$

because  $\bar{h}_X(\eta^*, \eta^*) = 0 \leq \varepsilon$ . On the other hand, if  $\eta \in A$ , choose  $j$  such that  $\bar{h}_X(\eta, \eta_j) \leq \varepsilon$ . Then

$$\mathbb{E}_\eta[1 - \phi_{A, \varepsilon, X} | X_{1:m}] \leq \mathbb{E}_\eta[1 - \phi_j | X_{1:m}] \leq e^{-m\varepsilon^2}.$$

This proves Equations (59) and (60).  $\square$

### D.3. Entropy bounds

We next control the metric entropy of the reduced experiment on the high-probability good-design event introduced in Appendix B.

For  $i \in [m]$ , let

$$\mathcal{E}_n^{(i)}(\delta) := \{(1 - \delta)nI_d \preceq X_i X_i^\top \preceq (1 + \delta)nI_d\},$$

and define the joint good-design event

$$\mathcal{G}_{m,n}(\delta) := \bigcap_{i=1}^m \mathcal{E}_n^{(i)}(\delta). \quad (61)$$

**Lemma D.4** (Probability of the joint good-design event). *For every  $0 < \delta < 1$ , there exist constants  $c_{d,\delta}, C_{d,\delta} \in (0, \infty)$ , depending only on  $d$  and  $\delta$ , such that*

$$\mathbb{P}_{\eta^*}(\mathcal{G}_{m,n}(\delta)^c) \leq m C_{d,\delta} e^{-c_{d,\delta} n}.$$

*Proof.* This follows immediately from Theorem B.10 and a union bound over  $i = 1, \dots, m$ .  $\square$

We will repeatedly use the following discrete Hellinger distance on the simplex:

$$h_{\text{disc}}^2(\pi, \pi') := \frac{1}{2} \sum_{z \in \mathcal{L}} (\sqrt{\pi_z} - \sqrt{\pi'_z})^2.$$

**Lemma D.5** (Componentwise telescoping bound for composed operators). *There exists a constant*

$$C_{\text{prod}} = C_{\text{prod}}(L, B, r) < \infty$$

*such that for every  $\Theta = (A_1, \dots, A_r) \in \mathbb{A}_B^r$ , every  $\Theta' = (A'_1, \dots, A'_r) \in \mathbb{A}_B^r$ , and every  $z \in \mathcal{L}$ ,*

$$\|A_z(\Theta) - A_z(\Theta')\|_F \leq C_{\text{prod}} \|\Theta - \Theta'\|_{\text{lib}}. \quad (62)$$

*One may take  $C_{\text{prod}} = LB^{L-1}\sqrt{r}$ .*

*Proof.* Fix  $z = (z_1, \dots, z_t) \in \mathcal{L}$ , with  $t \leq L$ . A telescoping expansion gives

$$A_z(\Theta) - A_z(\Theta') = \sum_{\ell=1}^t A_{z_\ell} \cdots A_{z_{\ell+1}} (A_{z_\ell} - A'_{z_\ell}) A'_{z_{\ell-1}} \cdots A'_{z_1}.$$

Taking Frobenius norms, using submultiplicativity of the operator norm, and the bound  $\|A_a\|_{\text{op}}, \|A'_a\|_{\text{op}} \leq B$ , we obtain

$$\|A_z(\Theta) - A_z(\Theta')\|_F \leq \sum_{\ell=1}^t B^{t-1} \|A_{z_\ell} - A'_{z_\ell}\|_F \leq LB^{L-1} \sum_{a=1}^r \|A_a - A'_a\|_F.$$

Finally,

$$\sum_{a=1}^r \|A_a - A'_a\|_F \leq \sqrt{r} \|\Theta - \Theta'\|_{\text{lib}}.$$

Combining the displays yields Equation (62).  $\square$

The next lemma is the key predictive upper bound. It shows that, on the good-design event, the reduced conditional Hellinger geometry is locally equivalent to the direct-product metric

$$(\pi, \Theta) \mapsto h_{\text{disc}}(\pi, \pi^*) + \sqrt{n} \|\Theta - \Theta^*\|_{\text{lib}}.$$

**Lemma D.6** (Local Hellinger upper bound in the reduced experiment). *Fix  $0 < \delta < 1$ . There exist constants  $c_{\text{loc}}, C_{\text{loc}} \in (0, \infty)$ , depending only on  $(\mathcal{L}, r, d, L, B, \sigma, \delta)$ , such that the following holds on  $\mathcal{G}_{m,n}(\delta)$ . If  $\eta = (\pi, \Theta)$  and  $\eta' = (\pi', \Theta')$  satisfy*

$$\|\Theta - \Theta'\|_{\text{lib}} \leq c_{\text{loc}} n^{-1/2},$$

then

$$\bar{h}_X^2(\eta, \eta') \leq 2h_{\text{disc}}^2(\pi, \pi') + C_{\text{loc}} n \|\Theta - \Theta'\|_{\text{lib}}^2. \quad (63)$$

*Proof.* Fix a task index  $i \in [m]$ , and write

$$f_{z,i} := \mathcal{MN}_{d,d}(A_z(\Theta), \sigma^2 I_d, (X_i X_i^\top)^{-1}), \quad f'_{z,i} := \mathcal{MN}_{d,d}(A_z(\Theta'), \sigma^2 I_d, (X_i X_i^\top)^{-1}).$$

Then

$$q_{\eta,i} = \sum_{z \in \mathcal{L}} \pi_z f_{z,i}, \quad q_{\eta',i} = \sum_{z \in \mathcal{L}} \pi'_z f'_{z,i}.$$

By the triangle inequality for Hellinger distance,

$$h(q_{\eta,i}, q_{\eta',i}) \leq h\left(\sum_z \pi_z f_{z,i}, \sum_z \pi_z f'_{z,i}\right) + h\left(\sum_z \pi_z f'_{z,i}, \sum_z \pi'_z f'_{z,i}\right).$$

Hence

$$h^2(q_{\eta,i}, q_{\eta',i}) \leq 2h_1^2 + 2h_2^2, \quad (64)$$

where

$$h_1 := h\left(\sum_z \pi_z f_{z,i}, \sum_z \pi_z f'_{z,i}\right), \quad h_2 := h\left(\sum_z \pi_z f'_{z,i}, \sum_z \pi'_z f'_{z,i}\right).$$

For the weights term, consider the joint laws on  $\mathcal{L} \times \mathbb{R}^{d \times d}$ ,

$$\tilde{q}_i(z, w) := \pi_z f'_{z,i}(w), \quad \tilde{q}'_i(z, w) := \pi'_z f'_{z,i}(w).$$

Marginalizing over  $z$  is a measurable map, so by data processing for Hellinger distance,

$$h_2^2 \leq h^2(\tilde{q}_i, \tilde{q}'_i) = h_{\text{disc}}^2(\pi, \pi').$$

For the component-location term, Hellinger affinity is concave in each argument, hence

$$\rho\left(\sum_z \pi_z f_{z,i}, \sum_z \pi_z f'_{z,i}\right) \geq \sum_{z \in \mathcal{L}} \pi_z \rho(f_{z,i}, f'_{z,i}).$$

Therefore

$$h_1^2 \leq \sum_{z \in \mathcal{L}} \pi_z h^2(f_{z,i}, f'_{z,i}).$$

By Theorem D.5,

$$\max_{z \in \mathcal{L}} \|A_z(\Theta) - A_z(\Theta')\|_F \leq C_{\text{prod}} \|\Theta - \Theta'\|_{\text{lib}}.$$

Hence, if  $\|\Theta - \Theta'\|_{\text{lib}} \leq c_{\text{loc}} n^{-1/2}$  and  $c_{\text{loc}}$  is chosen sufficiently small, then

$$\frac{(1 + \delta)n}{8\sigma^2} \max_{z \in \mathcal{L}} \|A_z(\Theta) - A_z(\Theta')\|_F^2 \leq 1.$$

On  $\mathcal{G}_{m,n}(\delta)$ , Theorem B.11 from Appendix B then yields

$$h^2(f_{z,i}, f'_{z,i}) \leq C'_{\text{loc}} n \|A_z(\Theta) - A_z(\Theta')\|_F^2 \leq C'_{\text{loc}} n C_{\text{prod}}^2 \|\Theta - \Theta'\|_{\text{lib}}^2.$$

Summing over  $z$  with weights  $\pi_z$  gives

$$h_1^2 \leq C'_{\text{loc}} C_{\text{prod}}^2 n \|\Theta - \Theta'\|_{\text{lib}}^2.$$

Substituting the bounds for  $h_1^2$  and  $h_2^2$  into Equation (64) yields

$$h^2(q_{\eta,i}, q_{\eta',i}) \leq 2h_{\text{disc}}^2(\pi, \pi') + 2C'_{\text{loc}} C_{\text{prod}}^2 n \|\Theta - \Theta'\|_{\text{lib}}^2.$$

Averaging over  $i = 1, \dots, m$  proves Equation (63).  $\square$

We now convert the preceding local Hellinger bound into a covering-number estimate.

**Proposition D.7** (Conditional metric entropy). *Fix  $0 < \delta < 1$ . There exist constants  $\varepsilon_0 \in (0, 1)$  and  $C_{\text{ent},1}, C_{\text{ent},2}, C_{\text{ent},3}, C_{\text{ent},4} \in (0, \infty)$ , depending only on  $(\mathcal{L}, r, d, L, B, \sigma, \delta)$ , such that on  $\mathcal{G}_{m,n}(\delta)$ , for every  $0 < \varepsilon \leq \varepsilon_0$ ,*

$$\log N(\varepsilon, \Xi(B, \mathcal{L}), \bar{h}_X) \leq C_{\text{ent},1} K_{\text{eff}} \log \frac{C_{\text{ent},2}}{\varepsilon} + C_{\text{ent},3} r d^2 \log \frac{C_{\text{ent},4} \sqrt{n}}{\varepsilon}. \quad (65)$$

The same bound holds, up to an additive constant  $\log |\text{Aut}(\mathcal{L})|$ , for the quotient space  $\Xi(B, \mathcal{L}) / \text{Aut}(\mathcal{L})$ .

*Proof.* We cover the simplex and the module library separately.

**Step 1: covering the grammar weights.** Since  $\mathcal{L}$  is finite and

$$\mathfrak{P}_{\mathcal{L}}^{\circ} \subseteq \left\{ \pi \in [0, 1]^{|\mathcal{L}|} : \sum_z \pi_z = 1 \right\},$$

the simplex has Euclidean dimension  $K_{\text{eff}} = |\mathcal{L}| - 1$ . A standard volume argument yields

$$\log N(u, \mathfrak{P}_{\mathcal{L}}^{\circ}, \|\cdot\|_1) \leq C_1 K_{\text{eff}} \log \frac{C_2}{u} \quad \text{for } 0 < u \leq 1.$$

Since  $h_{\text{disc}}^2(\pi, \pi') \leq \frac{1}{2} \|\pi - \pi'\|_1$ , it follows that

$$\log N\left(\frac{\varepsilon}{4}, \mathfrak{P}_{\mathcal{L}}^{\circ}, h_{\text{disc}}\right) \leq C_1 K_{\text{eff}} \log \frac{C_2}{\varepsilon}.$$

**Step 2: covering the module library.** The compact set  $\mathbb{A}_B^r$  is contained in a Euclidean ball of radius  $B\sqrt{rd}$  in the norm  $\|\cdot\|_{\text{lib}}$ , and has ambient dimension  $rd^2$ . Hence

$$\log N(v, \mathbb{A}_B^r, \|\cdot\|_{\text{lib}}) \leq C_3 r d^2 \log \frac{C_4}{v} \quad \text{for } 0 < v \leq 1.$$

**Step 3: combine the coverings through Theorem D.6.** Choose the module covering radius

$$v_\varepsilon := \min \left\{ \frac{c_{\text{loc}}}{\sqrt{n}}, \frac{\varepsilon}{4\sqrt{C_{\text{loc}}n}} \right\},$$

and the simplex covering radius

$$u_\varepsilon := \frac{\varepsilon}{4}.$$

For any  $\eta = (\pi, \Theta)$ , select centers  $\pi^\sharp$  and  $\Theta^\sharp$  from the two coverings such that

$$h_{\text{disc}}(\pi, \pi^\sharp) \leq u_\varepsilon, \quad \|\Theta - \Theta^\sharp\|_{\text{lib}} \leq v_\varepsilon.$$

Then on  $\mathcal{G}_{m,n}(\delta)$ , Theorem D.6 implies

$$\bar{h}_X^2((\pi, \Theta), (\pi^\sharp, \Theta^\sharp)) \leq 2u_\varepsilon^2 + C_{\text{loc}}nv_\varepsilon^2 \leq \frac{\varepsilon^2}{8} + \frac{\varepsilon^2}{16} < \varepsilon^2.$$

Thus the Cartesian product of the two coverings is an  $\varepsilon$ -cover of  $\Xi(B, \mathcal{L})$  under  $\bar{h}_X$ . Taking logarithms and absorbing constants yields Equation (65). Since passing to the quotient cannot increase covering numbers by more than the finite factor  $|\text{Aut}(\mathcal{L})|$ , the quotient statement follows.  $\square$

#### D.4. Local prior thickness

We now verify that the prior allocates enough mass to neighborhoods whose conditional Kullback–Leibler size is  $O(\varepsilon^2)$ . The key simplification is that KL for the observed mixture can be bounded by KL in the augmented experiment where the latent program  $Z$  is observed.

Fix the truth  $\eta^* = (\pi^*, \Theta^*)$ . For  $X \in \mathbb{R}^{d \times n}$  full row rank and  $\eta = (\pi, \Theta)$ , define the augmented one-task conditional density

$$\tilde{q}_{\eta^*, X}(z, w) := \pi_z \varphi_X(w; A_z(\Theta)), \quad z \in \mathcal{L}, w \in \mathbb{R}^{d \times d}, \quad (66)$$

where  $\varphi_X(\cdot; W)$  is the density of  $\mathcal{MN}_{d,d}(W, \sigma^2 I_d, (XX^\top)^{-1})$ .

**Lemma D.8** (KL upper bound via latent augmentation). *For every full-row-rank design  $X$  and every  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$ ,*

$$\text{KL}(q_{\eta^*}(\cdot | X) \| q_\eta(\cdot | X)) \leq \text{KL}(\pi^*, \pi) + \sum_{z \in \mathcal{L}} \pi_z^* \text{KL}(\mathcal{N}_X^{A_z^*} \| \mathcal{N}_X^{A_z}). \quad (67)$$

*Proof.* Consider the joint laws of  $(Z, \widehat{W})$  under  $\eta^*$  and  $\eta$ , namely

$$\tilde{q}_{\eta^*, X}(z, w) = \pi_z^* \varphi_X(w; A_z^*), \quad \tilde{q}_{\eta, X}(z, w) = \pi_z \varphi_X(w; A_z).$$

Marginalizing over  $z$  maps these joint laws to the observed laws  $q_{\eta^*}(\cdot | X)$  and  $q_\eta(\cdot | X)$ . Since Kullback–Leibler divergence contracts under measurable maps,

$$\text{KL}(q_{\eta^*}(\cdot | X) \| q_\eta(\cdot | X)) \leq \text{KL}(\tilde{q}_{\eta^*, X} \| \tilde{q}_{\eta, X}).$$

Because the joint density factorizes into a discrete law on  $z$  and a conditional Gaussian law on  $\widehat{W}$ ,

$$\text{KL}(\tilde{q}_{\eta^*, X} \| \tilde{q}_{\eta, X}) = \sum_{z \in \mathcal{L}} \pi_z^* \log \frac{\pi_z^*}{\pi_z} + \sum_{z \in \mathcal{L}} \pi_z^* \text{KL}(\mathcal{N}_X^{A_z^*} \| \mathcal{N}_X^{A_z}),$$

which is exactly Equation (67).  $\square$

To describe local prior thickness, define the neighborhood

$$\mathcal{U}_n(\varepsilon) := \left\{ (\pi, \Theta) \in \Xi(B, \mathcal{L}) : \|\pi - \pi^*\|_1 \leq c_\pi \varepsilon, \|\Theta - \Theta^*\|_{\text{lib}} \leq c_\Theta \frac{\varepsilon}{\sqrt{n}} \right\}, \quad (68)$$

where  $c_\pi, c_\Theta > 0$  will be chosen sufficiently small.

**Lemma D.9** (Local KL control on  $\mathcal{U}_n(\varepsilon)$ ). *Fix  $0 < \delta < 1$ . There exist constants  $c_\pi, c_\Theta, \varepsilon_0, C_{\text{KL}} \in (0, \infty)$ , depending only on  $(\mathcal{L}, r, d, L, B, \sigma, \delta, \pi_{\min})$ , such that on  $\mathcal{G}_{m,n}(\delta)$ , for every  $0 < \varepsilon \leq \varepsilon_0$  and every  $\eta \in \mathcal{U}_n(\varepsilon)$ ,*

$$\frac{1}{m} \sum_{i=1}^m \text{KL}(q_{\eta^*}(\cdot | X_i) \| q_\eta(\cdot | X_i)) \leq C_{\text{KL}} \varepsilon^2. \quad (69)$$

*Proof.* Fix  $\eta = (\pi, \Theta) \in \mathcal{U}_n(\varepsilon)$ . Because  $\pi_z^* \geq \pi_{\min}$  for all  $z \in \mathcal{L}$  and  $\|\pi - \pi^*\|_1 \leq c_\pi \varepsilon$ , choosing  $c_\pi \varepsilon_0 \leq \pi_{\min}/2$  ensures

$$\pi_z \geq \frac{\pi_{\min}}{2} \quad \text{for all } z \in \mathcal{L}.$$

Using the elementary inequality  $\log x \leq x - 1$ , one obtains

$$\text{KL}(\pi^*, \pi) = \sum_{z \in \mathcal{L}} \pi_z^* \log \frac{\pi_z^*}{\pi_z} \leq \sum_{z \in \mathcal{L}} \frac{(\pi_z^* - \pi_z)^2}{\pi_z} \leq \frac{2}{\pi_{\min}} \|\pi - \pi^*\|_2^2 \leq \frac{2}{\pi_{\min}} \|\pi - \pi^*\|_1^2 \leq \frac{2c_\pi^2}{\pi_{\min}} \varepsilon^2.$$

Next, by Theorem D.5,

$$\max_{z \in \mathcal{L}} \|A_z - A_z^*\|_F \leq C_{\text{prod}} \|\Theta - \Theta^*\|_{\text{lib}} \leq C_{\text{prod}} c_\Theta \frac{\varepsilon}{\sqrt{n}}.$$

On  $\mathcal{G}_{m,n}(\delta)$ , Equation (37) from Appendix B gives, for each  $i \in [m]$  and each  $z \in \mathcal{L}$ ,

$$\text{KL}(\mathcal{N}_{X_i}^{A_z^*} \| \mathcal{N}_{X_i}^{A_z}) \leq \frac{(1 + \delta)n}{2\sigma^2} \|A_z^* - A_z\|_F^2 \leq \frac{(1 + \delta)C_{\text{prod}}^2 c_\Theta^2}{2\sigma^2} \varepsilon^2.$$

Applying Theorem D.8 yields

$$\text{KL}(q_{\eta^*}(\cdot | X_i) \| q_\eta(\cdot | X_i)) \leq \frac{2c_\pi^2}{\pi_{\min}} \varepsilon^2 + \frac{(1 + \delta)C_{\text{prod}}^2 c_\Theta^2}{2\sigma^2} \varepsilon^2.$$

The right-hand side is independent of  $i$ , so averaging over  $i = 1, \dots, m$  gives Equation (69).  $\square$

We now lower-bound the prior mass of  $\mathcal{U}_n(\varepsilon)$ .

**Proposition D.10** (Local prior thickness). *There exist constants  $\varepsilon_0 \in (0, 1)$  and  $c_\Pi, C_\Pi \in (0, \infty)$ , depending only on  $(\mathcal{L}, r, d, L, B, \sigma, \Pi, \eta^*)$ , such that for every  $0 < \varepsilon \leq \varepsilon_0$ ,*

$$\Pi(\mathcal{U}_n(\varepsilon)) \geq c_\Pi \varepsilon^{K_{\text{eff}} + rd^2} n^{-rd^2/2}, \quad (70)$$

and therefore

$$-\log \Pi(\mathcal{U}_n(\varepsilon)) \leq C_\Pi \left[ K_{\text{eff}} \log \frac{1}{\varepsilon} + rd^2 \log \frac{\sqrt{n}}{\varepsilon} \right]. \quad (71)$$

*Proof.* Because  $\pi^*$  lies in the interior of the simplex and  $\Pi_\pi$  has a density that is continuous and strictly positive near  $\pi^*$ , there exist  $r_\pi > 0$  and  $c'_\pi > 0$  such that the prior density of  $\Pi_\pi$  is bounded below by  $c'_\pi$  on the  $\ell_1$ -ball  $B_1(\pi^*, r_\pi) \cap \mathfrak{P}_\mathcal{L}^\circ$ . Since the simplex has dimension  $K_{\text{eff}} = |\mathcal{L}| - 1$ , the  $K_{\text{eff}}$ -dimensional volume of

$$\{\pi \in \mathfrak{P}_\mathcal{L}^\circ : \|\pi - \pi^*\|_1 \leq c_\pi \varepsilon\}$$

is bounded below by  $c''_\pi \varepsilon^{K_{\text{eff}}}$  for sufficiently small  $\varepsilon$ . Hence

$$\Pi_\pi(\|\pi - \pi^*\|_1 \leq c_\pi \varepsilon) \geq c'_\pi c''_\pi \varepsilon^{K_{\text{eff}}}.$$

Similarly, because  $\Pi_\Theta = \bigotimes_{a=1}^r \Pi_a$  has a continuous strictly positive density near  $\Theta^*$ , there exist  $r_\Theta > 0$  and  $c'_\Theta > 0$  such that this density is bounded below by  $c'_\Theta$  on the  $\|\cdot\|_{\text{lib}}$ -ball  $B_{\text{lib}}(\Theta^*, r_\Theta)$ . The ambient dimension of  $\mathbb{A}_B^r$  is  $rd^2$ , so the Euclidean volume of

$$\left\{ \Theta \in \mathbb{A}_B^r : \|\Theta - \Theta^*\|_{\text{lib}} \leq c_\Theta \frac{\varepsilon}{\sqrt{n}} \right\}$$

is bounded below by

$$c''_{\Theta} \left( \frac{\varepsilon}{\sqrt{n}} \right)^{rd^2}$$

for sufficiently small  $\varepsilon$ . Therefore

$$\Pi_{\Theta} \left( \|\Theta - \Theta^*\|_{\text{lib}} \leq c_{\Theta} \frac{\varepsilon}{\sqrt{n}} \right) \geq c'_{\Theta} c''_{\Theta} \left( \frac{\varepsilon}{\sqrt{n}} \right)^{rd^2}.$$

Using the product structure  $\Pi = \Pi_{\pi} \otimes \Pi_{\Theta}$ , we conclude that

$$\Pi(\mathcal{U}_n(\varepsilon)) \geq (c'_{\pi} c''_{\pi})(c'_{\Theta} c''_{\Theta}) \varepsilon^{K_{\text{eff}} + rd^2} n^{-rd^2/2},$$

which is Equation (70). Taking negative logarithms yields Equation (71) after enlarging the constant.  $\square$

### D.5. Packaging the contraction ingredients

The final proposition records the exact balance needed in Appendix E.

**Proposition D.11** (Entropy and prior mass at the predictive radius). *There exists  $M_0 \in (0, \infty)$ , depending only on the fixed model class and the prior, such that the following holds. Fix  $0 < \delta < 1$ , and let  $\delta_{m,n}$  be as in Equation (57). For every  $M \geq M_0$ , if*

$$\varepsilon_{m,n} := M \delta_{m,n},$$

then for all sufficiently large  $m$  and all  $n \geq d + 1$ , on the event  $\mathcal{G}_{m,n}(\delta)$ ,

$$\log N(\varepsilon_{m,n}, \Xi(B, \mathcal{L}), \bar{h}_X) \leq \frac{m\varepsilon_{m,n}^2}{8}, \quad (72)$$

$$-\log \Pi(\mathcal{U}_n(\varepsilon_{m,n})) \leq \frac{m\varepsilon_{m,n}^2}{8}. \quad (73)$$

*Proof.* By Theorem D.7,

$$\log N(\varepsilon, \Xi(B, \mathcal{L}), \bar{h}_X) \leq C_1 K_{\text{eff}} \log \frac{C_2}{\varepsilon} + C_3 rd^2 \log \frac{C_4 \sqrt{n}}{\varepsilon}.$$

Likewise, Theorem D.10 gives

$$-\log \Pi(\mathcal{U}_n(\varepsilon)) \leq C_5 K_{\text{eff}} \log \frac{1}{\varepsilon} + C_5 rd^2 \log \frac{\sqrt{n}}{\varepsilon}.$$

Substitute  $\varepsilon = \varepsilon_{m,n} = M \delta_{m,n}$ , with

$$m \delta_{m,n}^2 = K_{\text{eff}} \log m + rd^2 \log(mn).$$

Since  $\delta_{m,n} \rightarrow 0$  along any asymptotic regime with  $m \rightarrow \infty$  and  $n \geq d + 1$ , the logarithmic terms satisfy

$$\log \frac{1}{\varepsilon_{m,n}} \lesssim \log(mn)$$

for all sufficiently large  $m$ . Hence both the entropy bound and the prior-thickness bound are controlled by

$$C [K_{\text{eff}} \log m + rd^2 \log(mn)] = C m \delta_{m,n}^2.$$

Choosing  $M_0$  large enough so that  $C \leq M_0^2/8$  yields Equations (72) and (73) for every  $M \geq M_0$ .  $\square$

*Remark D.12* (What this appendix has established). Appendix D proves the full contraction toolkit for the reduced experiment:

1. exponentially consistent conditional tests for  $\bar{h}_X$ -separated alternatives;
2. conditional entropy bounds at the predictive scale  $\delta_{m,n}$ ;
3. local prior thickness in Kullback–Leibler neighborhoods of size  $O(\varepsilon^2)$ .

In Appendix E, these ingredients are combined with the exact reduction of Appendix B to obtain predictive posterior contraction for the grouped task-family law.

## E. Predictive Posterior Contraction and Latent-Program Recovery

This appendix combines the exact reduction of Appendix B with the testing / entropy / prior-thickness bounds of Appendix D. The resulting predictive contraction theorem is proved in the natural conditional metric  $\bar{h}_X$ , i.e., the average Hellinger distance in the reduced non-i.i.d. experiment given the random designs. This is the correct predictive metric for the grouped task-family law after sufficient-statistic reduction. The proof follows the standard testing-based posterior-contraction program for independent non-identically distributed observations (Ghosal & van der Vaart, 2007; Ghosal et al., 2000), but every model-specific ingredient has already been established directly in the preceding appendices.

### E.1. Posterior notation and design regularity regime

For the reduced experiment, define the posterior

$$\Pi_m(B) := \Pi(B \mid X_{1:m}, \widehat{W}_{1:m}) = \frac{\int_B \prod_{i=1}^m q_\eta(\widehat{W}_i \mid X_i) \Pi(d\eta)}{\int_{\Xi(B, \mathcal{L})} \prod_{i=1}^m q_\eta(\widehat{W}_i \mid X_i) \Pi(d\eta)}, \quad B \subseteq \Xi(B, \mathcal{L}) \text{ measurable.} \quad (74)$$

By Theorem B.7 from Appendix B, this equals the posterior based on the full grouped data  $(X_i, Y_i)_{i=1}^m$ .

Throughout this appendix we work under the asymptotic regime

$$m e^{-c_{d,\delta} n} \rightarrow 0 \quad \text{for some fixed } 0 < \delta < 1, \quad (75)$$

where  $c_{d,\delta}$  is the constant from Theorem D.4. Under Equation (75), the joint good-design event

$$\mathcal{G}_{m,n}(\delta) = \bigcap_{i=1}^m \{(1 - \delta)nI_d \preceq X_i X_i^\top \preceq (1 + \delta)nI_d\}$$

has probability tending to one under the true law.

We also recall the predictive radius from Appendix D,

$$\delta_{m,n}^2 = \frac{K_{\text{eff}} \log m + r d^2 \log(mn)}{m}.$$

### E.2. Predictive posterior contraction in the reduced experiment

The contraction theorem is most naturally stated in the empirical conditional Hellinger metric

$$\bar{h}_X^2(\eta, \eta^*) = \frac{1}{m} \sum_{i=1}^m h^2(q_\eta(\cdot \mid X_i), q_{\eta^*}(\cdot \mid X_i)).$$

**Theorem E.1** (Predictive posterior contraction in  $\bar{h}_X$ ). *Assume Theorems A.1 and A.2, and suppose that Equation (75) holds for some fixed  $0 < \delta < 1$ . Then there exists a constant  $M_{\text{pred}} \in (0, \infty)$ , depending only on the fixed model class and the prior, such that for every  $M \geq M_{\text{pred}}$ ,*

$$\Pi_m(\eta \in \Xi(B, \mathcal{L}) : \bar{h}_X(\eta, \eta^*) > M \delta_{m,n}) \rightarrow 0 \quad (76)$$

in  $P_{\eta^*}$ -probability.

*Proof.* Fix  $0 < \delta < 1$ , and let  $\mathcal{G}_{m,n}(\delta)$  be the joint good-design event from Equation (61). By Theorem D.4 and the regime Equation (75),

$$P_{\eta^*}(\mathcal{G}_{m,n}(\delta)^c) \rightarrow 0.$$

It therefore suffices to prove Equation (76) on  $\mathcal{G}_{m,n}(\delta)$ .

Conditional on  $X_{1:m}$ , the reduced observations  $\widehat{W}_1, \dots, \widehat{W}_m$  are independent, though generally non-identically distributed, with conditional densities  $q_{\eta,1}, \dots, q_{\eta,m}$ . Hence the general testing-based contraction theorem for independent non-i.i.d. experiments from Ghosal & van der Vaart (2007) applies once its three hypotheses are verified for the semimetric

$$d_m(\eta, \eta') := \bar{h}_X(\eta, \eta').$$

We now check those hypotheses on  $\mathcal{G}_{m,n}(\delta)$ .

**Testing condition.** For any  $\varepsilon > 0$  and any measurable set  $A \subseteq \Xi(B, \mathcal{L})$  satisfying

$$\inf_{\eta \in A} \bar{h}_X(\eta, \eta^*) \geq 4\varepsilon,$$

Theorem D.3 constructs a measurable test  $\phi_{A, \varepsilon, X}$  with type-I and uniform type-II error bounds of the form

$$\mathbb{E}_{\eta^*}[\phi_{A, \varepsilon, X} \mid X_{1:m}] \leq N(\varepsilon, A, \bar{h}_X) e^{-m\varepsilon^2}, \quad \sup_{\eta \in A} \mathbb{E}_{\eta}[1 - \phi_{A, \varepsilon, X} \mid X_{1:m}] \leq e^{-m\varepsilon^2}.$$

Thus the testing hypothesis holds in the metric  $d_m = \bar{h}_X$ .

**Entropy condition.** By Theorem D.7, for all sufficiently small  $\varepsilon > 0$ ,

$$\log N(\varepsilon, \Xi(B, \mathcal{L}), \bar{h}_X) \leq C_{\text{ent},1} K_{\text{eff}} \log \frac{C_{\text{ent},2}}{\varepsilon} + C_{\text{ent},3} r d^2 \log \frac{C_{\text{ent},4} \sqrt{n}}{\varepsilon}.$$

Choosing  $\varepsilon = \varepsilon_{m,n} := M\delta_{m,n}$  and  $M \geq M_0$ , where  $M_0$  is the constant from Theorem D.11, yields

$$\log N(\varepsilon_{m,n}, \Xi(B, \mathcal{L}), \bar{h}_X) \leq \frac{m\varepsilon_{m,n}^2}{8}.$$

The same bound obviously controls the entropy of each shell

$$\{\eta : j\varepsilon_{m,n} < \bar{h}_X(\eta, \eta^*) \leq 2j\varepsilon_{m,n}\}, \quad j \geq 1,$$

since each shell is a subset of the full parameter space.

**Prior mass of KL neighborhoods.** By Theorem D.9, every  $\eta \in \mathcal{U}_n(\varepsilon)$  satisfies

$$\frac{1}{m} \sum_{i=1}^m \text{KL}(q_{\eta^*}(\cdot \mid X_i) \parallel q_{\eta}(\cdot \mid X_i)) \leq C_{\text{KL}} \varepsilon^2$$

on  $\mathcal{G}_{m,n}(\delta)$ . By Theorems D.10 and D.11, for  $\varepsilon = \varepsilon_{m,n} = M\delta_{m,n}$  and  $M$  sufficiently large,

$$-\log \Pi(\mathcal{U}_n(\varepsilon_{m,n})) \leq \frac{m\varepsilon_{m,n}^2}{8}.$$

Hence the prior-thickness hypothesis of the general theorem is satisfied.

We have therefore verified all conditions of the main posterior-contraction theorem of Ghosal & van der Vaart (2007) for the semimetric  $d_m = \bar{h}_X$  and the rate  $\varepsilon_{m,n} = M\delta_{m,n}$ . It follows that

$$\Pi_m(\eta : \bar{h}_X(\eta, \eta^*) > M\delta_{m,n}) \rightarrow 0$$

in conditional  $P_{\eta^*}(\cdot \mid X_{1:m})$ -probability, for every design sequence in  $\mathcal{G}_{m,n}(\delta)$ . Since  $P_{\eta^*}(\mathcal{G}_{m,n}(\delta)^c) \rightarrow 0$ , the same conclusion holds unconditionally in  $P_{\eta^*}$ -probability.  $\square$

**Corollary E.2** (Predictive contraction for the full grouped-data posterior). *Under the assumptions of Theorem E.1, the full-data posterior satisfies*

$$\Pi(\eta \in \Xi(B, \mathcal{L}) : \bar{h}_X(\eta, \eta^*) > M\delta_{m,n} \mid X_{1:m}, Y_{1:m}) \longrightarrow 0 \quad (77)$$

in  $P_{\eta^*}$ -probability for every  $M \geq M_{\text{pred}}$ .

*Proof.* This is an immediate consequence of Theorem E.1 and the exact posterior-reduction identity Equation (31).  $\square$

### E.3. The reduced predictive law of one task

The next lemma records the precise predictive object estimated by the metric  $\bar{h}_X$ .

**Definition E.3** (Reduced one-task predictive law). For  $\eta \in \Xi(B, \mathcal{L})$ , define the reduced one-task predictive law

$$\tilde{Q}_\eta^{(n)}(dX, dW) := p_X(X) q_\eta(W | X) dX dW,$$

where  $p_X$  is the Gaussian design density and  $q_\eta(\cdot | X)$  is the conditional density of  $\widehat{W}$  given  $X$ .

**Lemma E.4** (Exact Hellinger representation for the reduced predictive law). For every  $\eta, \eta' \in \Xi(B, \mathcal{L})$ ,

$$h^2(\tilde{Q}_\eta^{(n)}, \tilde{Q}_{\eta'}^{(n)}) = \mathbb{E}_X [h^2(q_\eta(\cdot | X), q_{\eta'}(\cdot | X))]. \quad (78)$$

Consequently,

$$\mathbb{E}[\bar{h}_X^2(\eta, \eta')] = h^2(\tilde{Q}_\eta^{(n)}, \tilde{Q}_{\eta'}^{(n)}).$$

*Proof.* By definition,

$$\rho(\tilde{Q}_\eta^{(n)}, \tilde{Q}_{\eta'}^{(n)}) = \int p_X(x) \left[ \int \sqrt{q_\eta(w | x) q_{\eta'}(w | x)} dw \right] dx = \mathbb{E}_X [\rho(q_\eta(\cdot | X), q_{\eta'}(\cdot | X))].$$

Since  $h^2 = 1 - \rho$ , we obtain

$$h^2(\tilde{Q}_\eta^{(n)}, \tilde{Q}_{\eta'}^{(n)}) = 1 - \mathbb{E}_X [\rho(q_\eta, q_{\eta'})] = \mathbb{E}_X [1 - \rho(q_\eta, q_{\eta'})] = \mathbb{E}_X [h^2(q_\eta(\cdot | X), q_{\eta'}(\cdot | X))].$$

The second claim follows because  $X_1, \dots, X_m$  are i.i.d. □

*Remark E.5* (Interpretation of the predictive theorem). Theorem E.1 is therefore a posterior contraction theorem for the reduced one-task predictive law  $\tilde{Q}_\eta^{(n)}$ , measured through its empirical conditional Hellinger analogue  $\bar{h}_X$ . This is the natural predictive object after exact sufficient-statistic reduction.

### E.4. Oracle latent-program recovery

We now quantify how quickly the latent task program can be recovered once the true task-family parameter is known. This is the cleanest posterior decoding statement available before the structural-localization theory of Appendix F.

For  $i \in [m]$ ,  $z \in \mathcal{L}$ , and  $\eta = (\pi, \Theta)$ , define the component responsibility

$$\omega_{i,z}(\eta) := \Pi_\eta(Z_i = z | X_i, \widehat{W}_i) = \frac{\pi_z \varphi_{X_i}(\widehat{W}_i; A_z(\Theta))}{\sum_{u \in \mathcal{L}} \pi_u \varphi_{X_i}(\widehat{W}_i; A_u(\Theta))}, \quad (79)$$

where  $\varphi_{X_i}(\cdot; A_z(\Theta))$  denotes the density of  $\mathcal{MN}_{d,d}(A_z(\Theta), \sigma^2 I_d, (X_i X_i^\top)^{-1})$ .

Recall the true component-separation margin from Equation (53):

$$\gamma_{\text{comp}}(\eta^*) = \min_{z \neq z'} \|A_z^* - A_{z'}^*\|_F.$$

**Proposition E.6** (Oracle latent-program recovery under component separation). Assume Theorem A.1 and suppose that

$$\gamma_{\text{comp}}(\eta^*) > 0.$$

Fix  $0 < \delta < 1$ . Then on the joint good-design event  $\mathcal{G}_{m,n}(\delta)$ , for every  $i \in [m]$ ,

$$\mathbb{E}_{\eta^*} [1 - \omega_{i,Z_i^*}(\eta^*) | X_i, Z_i^*] \leq C_{\text{lat}} \exp\left(-\frac{(1-\delta)n \gamma_{\text{comp}}(\eta^*)^2}{8\sigma^2}\right), \quad (80)$$

where one may take

$$C_{\text{lat}} = \frac{|\mathcal{L}| - 1}{\sqrt{\pi_{\min}}}.$$

Consequently,

$$\frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\eta^*} [1 - \omega_{i,Z_i^*}(\eta^*) | X_{1:m}, Z_{1:m}^*] \leq C_{\text{lat}} \exp\left(-\frac{(1-\delta)n \gamma_{\text{comp}}(\eta^*)^2}{8\sigma^2}\right). \quad (81)$$

*Proof.* Fix a task index  $i$ , and condition on  $X_i$  and  $Z_i^* = z$ . Under the true parameter  $\eta^* = (\pi^*, \Theta^*)$ , the reduced observation  $\widehat{W}_i$  has density

$$f_{z,i}^*(w) = \varphi_{X_i}(w; A_z^*).$$

For  $u \neq z$ , let  $f_{u,i}^*(w) := \varphi_{X_i}(w; A_u^*)$ . Then

$$1 - \omega_{i,z}(\eta^*) = \sum_{u \neq z} \frac{\pi_u^* f_{u,i}^*(\widehat{W}_i)}{\sum_{v \in \mathcal{L}} \pi_v^* f_{v,i}^*(\widehat{W}_i)}.$$

For each  $u \neq z$ , set

$$t_u(\widehat{W}_i) := \frac{\pi_u^* f_{u,i}^*(\widehat{W}_i)}{\pi_z^* f_{z,i}^*(\widehat{W}_i)}.$$

Since  $\sum_v \pi_v^* f_{v,i}^* \geq \pi_z^* f_{z,i}^* + \pi_u^* f_{u,i}^*$ , we have

$$\frac{\pi_u^* f_{u,i}^*}{\sum_v \pi_v^* f_{v,i}^*} \leq \frac{t_u}{1 + t_u}.$$

Using the elementary inequality  $t/(1+t) \leq \sqrt{t}$  for  $t \geq 0$ , we obtain

$$\frac{\pi_u^* f_{u,i}^*}{\sum_v \pi_v^* f_{v,i}^*} \leq \sqrt{\frac{\pi_u^*}{\pi_z^*}} \sqrt{\frac{f_{u,i}^*}{f_{z,i}^*}}.$$

Taking expectation under  $f_{z,i}^*$  yields

$$\mathbb{E}_{\eta^*} \left[ \frac{\pi_u^* f_{u,i}^*(\widehat{W}_i)}{\sum_v \pi_v^* f_{v,i}^*(\widehat{W}_i)} \middle| X_i, Z_i^* = z \right] \leq \sqrt{\frac{\pi_u^*}{\pi_z^*}} \int \sqrt{f_{u,i}^*(w) f_{z,i}^*(w)} dw = \sqrt{\frac{\pi_u^*}{\pi_z^*}} \rho(f_{u,i}^*, f_{z,i}^*).$$

Summing over  $u \neq z$  gives

$$\mathbb{E}_{\eta^*} [1 - \omega_{i,z}(\eta^*) \mid X_i, Z_i^* = z] \leq \sum_{u \neq z} \sqrt{\frac{\pi_u^*}{\pi_z^*}} \rho(f_{u,i}^*, f_{z,i}^*). \quad (82)$$

Now on  $\mathcal{G}_{m,n}(\delta)$ , Equation (34) from Appendix B implies

$$\rho(f_{u,i}^*, f_{z,i}^*) = \exp\left(-\frac{1}{8\sigma^2} \|(A_u^* - A_z^*)X_i\|_F^2\right) \leq \exp\left(-\frac{(1-\delta)n}{8\sigma^2} \|A_u^* - A_z^*\|_F^2\right).$$

By the definition of  $\gamma_{\text{comp}}(\eta^*)$ ,

$$\|A_u^* - A_z^*\|_F \geq \gamma_{\text{comp}}(\eta^*) \quad \text{for every } u \neq z.$$

Hence

$$\rho(f_{u,i}^*, f_{z,i}^*) \leq \exp\left(-\frac{(1-\delta)n \gamma_{\text{comp}}(\eta^*)^2}{8\sigma^2}\right).$$

Also,  $\pi_z^* \geq \pi_{\min}$  and  $\pi_u^* \leq 1$ , so

$$\sqrt{\frac{\pi_u^*}{\pi_z^*}} \leq \pi_{\min}^{-1/2}.$$

Substituting these bounds into Equation (82) gives

$$\mathbb{E}_{\eta^*} [1 - \omega_{i,z}(\eta^*) \mid X_i, Z_i^* = z] \leq \frac{|\mathcal{L}| - 1}{\sqrt{\pi_{\min}}} \exp\left(-\frac{(1-\delta)n \gamma_{\text{comp}}(\eta^*)^2}{8\sigma^2}\right).$$

Since this bound is uniform in  $z \in \mathcal{L}$ , Equation (80) follows. Averaging over  $i = 1, \dots, m$  yields Equation (81).  $\square$

*Remark E.7* (From oracle recovery to fully Bayesian recovery). Theorem E.6 quantifies the intrinsic within-task difficulty of decoding the latent program when the true task-family parameter is known. To pass from the oracle decoder  $\omega_{i,z}(\eta^*)$  to the full hierarchical posterior decoder

$$\Pi\left(Z_i = z \mid X_{1:m}, \widehat{W}_{1:m}\right) = \int \omega_{i,z}(\eta) \Pi_m(d\eta),$$

one needs posterior localization in a genuine parameter neighborhood of  $\eta^*$ , not merely predictive contraction in  $\bar{h}_X$ . That structural localization is supplied later in Appendix F. Once combined with local continuity of  $\eta \mapsto \omega_{i,z}(\eta)$ , it yields the fully Bayesian latent-program recovery statement appearing in the main text.

## F. Local Inverse Inequalities and Structural Contraction

Singular statistical models exhibit a sharp distinction between predictive accuracy and parameter localization. This phenomenon is central in singular learning theory and in the weak-identifiability analysis of finite mixtures (Watanabe, 2009; Ho & Nguyen, 2016; 2019). The purpose of this appendix is to isolate the corresponding geometry for stochastic task grammars.

The appendix has three parts. First, we introduce a finite family of deconvolved Fourier witnesses that convert predictive Hellinger control in the reduced experiment into control of a finite-dimensional feature map of the latent task-family parameter. Second, under a regularity condition, we prove a local inverse inequality of order 1 and derive regular structural contraction of the quotient parameter. Third, we construct an explicit duplicated-module family in which the first-order geometry cancels, compute its reduced divergence exactly, and show that the local inverse exponent becomes 2.

### F.1. Finite deconvolved Fourier witnesses

Let

$$p := d^2.$$

For  $\eta = (\pi, \Theta) \in \Xi(B, \mathcal{L})$ , write

$$v_z(\Theta) := \text{vec}(A_z(\Theta)) \in \mathbb{R}^p, \quad \mu_\eta = \sum_{z \in \mathcal{L}} \pi_z \delta_{v_z(\Theta)}$$

for the induced mixing measure on operator space.

**Definition F.1** (Characteristic map of the induced mixing measure). For  $t \in \mathbb{R}^p$ , define

$$\psi_t(\eta) := \int e^{i\langle t, v \rangle} \mu_\eta(dv) = \sum_{z \in \mathcal{L}} \pi_z e^{i\langle t, v_z(\Theta) \rangle}. \quad (83)$$

For a finite set  $T = \{t_1, \dots, t_J\} \subset \mathbb{R}^p$ , define the finite witness map

$$\mathbb{T}_T(\eta) := (\psi_{t_1}(\eta), \dots, \psi_{t_J}(\eta)) \in \mathbb{C}^J. \quad (84)$$

**Lemma F.2** (Quotient invariance of the witness map). For every finite  $T \subset \mathbb{R}^p$ , the map  $\mathbb{T}_T$  is invariant under  $\text{Aut}(\mathcal{L})$ . Equivalently,

$$\mathbb{T}_T(\sigma \cdot \eta) = \mathbb{T}_T(\eta) \quad \text{for all } \eta \in \Xi(B, \mathcal{L}), \sigma \in \text{Aut}(\mathcal{L}).$$

*Proof.* By Theorem A.4 and the definition of the induced mixing measure, the quotient action preserves  $\mu_\eta$ . Since  $\psi_t(\eta)$  is the characteristic function of  $\mu_\eta$  evaluated at  $t$ , each coordinate of  $\mathbb{T}_T$  is quotient-invariant.  $\square$

Fix deterministic designs  $X_{1:m}$ , and for each task define

$$\Sigma_i := (X_i X_i^\top)^{-1} \otimes \sigma^2 I_d \in \mathbb{R}^{p \times p}.$$

**Lemma F.3** (Deconvolved Fourier identity). *Fix  $t \in \mathbb{R}^p$ . For each task  $i \in [m]$ , define the complex-valued function*

$$g_{t,i}(W) := \exp\left(i\langle t, \text{vec}(W) \rangle + \frac{1}{2}t^\top \Sigma_i t\right), \quad W \in \mathbb{R}^{d \times d}. \quad (85)$$

Then, for every  $\eta \in \Xi(B, \mathcal{L})$ ,

$$\int g_{t,i}(W) q_\eta(W | X_i) dW = \psi_t(\eta). \quad (86)$$

*Proof.* By Theorem B.6 in Appendix B,

$$q_\eta(\cdot | X_i) = \sum_{z \in \mathcal{L}} \pi_z \mathcal{MN}_{d,d}(A_z(\Theta), \sigma^2 I_d, (X_i X_i^\top)^{-1}).$$

Fix  $z \in \mathcal{L}$ . If  $W \sim \mathcal{MN}_{d,d}(A_z(\Theta), \sigma^2 I_d, (X_i X_i^\top)^{-1})$ , then

$$\text{vec}(W) \sim \mathcal{N}(v_z(\Theta), \Sigma_i).$$

Hence

$$\mathbb{E}\left[e^{i\langle t, \text{vec}(W) \rangle}\right] = \exp\left(i\langle t, v_z(\Theta) \rangle - \frac{1}{2}t^\top \Sigma_i t\right).$$

Multiplying by  $e^{\frac{1}{2}t^\top \Sigma_i t}$  gives

$$\mathbb{E}[g_{t,i}(W)] = e^{i\langle t, v_z(\Theta) \rangle}.$$

Averaging with respect to the mixing weights  $\pi_z$  proves Equation (86).  $\square$

**Lemma F.4** (Witness stability under the predictive metric). *Fix  $0 < \delta < 1$ , a finite set  $T = \{t_1, \dots, t_J\} \subset \mathbb{R}^p$ , and the joint good-design event  $\mathcal{G}_{m,n}(\delta)$ . Then on  $\mathcal{G}_{m,n}(\delta)$ ,*

$$\|\mathbb{T}_T(\eta) - \mathbb{T}_T(\eta')\|_2 \leq C_T(\delta) \bar{h}_X(\eta, \eta') \quad \text{for all } \eta, \eta' \in \Xi(B, \mathcal{L}), \quad (87)$$

where one may take

$$C_T(\delta) := 2\sqrt{2J} \max_{1 \leq j \leq J} \exp\left(\frac{\sigma^2 \|t_j\|_2^2}{2(1-\delta)(d+1)}\right). \quad (88)$$

*Proof.* Fix  $t \in T$ . By Theorem F.3,

$$\psi_t(\eta) - \psi_t(\eta') = \frac{1}{m} \sum_{i=1}^m \int g_{t,i}(W) (q_\eta(W | X_i) - q_{\eta'}(W | X_i)) dW.$$

Hence

$$|\psi_t(\eta) - \psi_t(\eta')| \leq \frac{1}{m} \sum_{i=1}^m 2\|g_{t,i}\|_\infty \text{TV}(q_\eta(\cdot | X_i), q_{\eta'}(\cdot | X_i)).$$

On  $\mathcal{G}_{m,n}(\delta)$ ,

$$X_i X_i^\top \succeq (1-\delta)nI_d \succeq (1-\delta)(d+1)I_d,$$

so

$$\Sigma_i = (X_i X_i^\top)^{-1} \otimes \sigma^2 I_d \preceq \frac{\sigma^2}{(1-\delta)(d+1)} I_p.$$

Therefore

$$\|g_{t,i}\|_\infty \leq \exp\left(\frac{1}{2}t^\top \Sigma_i t\right) \leq \exp\left(\frac{\sigma^2 \|t\|_2^2}{2(1-\delta)(d+1)}\right) =: B_t(\delta).$$

Using  $\text{TV}(P, Q) \leq \sqrt{2}h(P, Q)$  and Cauchy–Schwarz,

$$\frac{1}{m} \sum_{i=1}^m \text{TV}(q_{\eta,i}, q_{\eta',i}) \leq \sqrt{2} \frac{1}{m} \sum_{i=1}^m h(q_{\eta,i}, q_{\eta',i}) \leq \sqrt{2} \bar{h}_X(\eta, \eta').$$

Hence

$$|\psi_t(\eta) - \psi_t(\eta')| \leq 2\sqrt{2} B_t(\delta) \bar{h}_X(\eta, \eta').$$

Summing over  $t_1, \dots, t_J$  in Euclidean norm gives Equation (87).  $\square$

## F.2. Regular local inverse inequality

We now prove a local inverse theorem of order 1 under a finite-dimensional regularity condition.

**Definition F.5** (Tangent space and tangent norm). At the truth  $\eta^* = (\pi^*, \Theta^*)$ , define the tangent space

$$\mathcal{T}_{\eta^*} := \left\{ (\alpha, H) : \alpha = (\alpha_z)_{z \in \mathcal{L}} \in \mathbb{R}^{\mathcal{L}}, \sum_{z \in \mathcal{L}} \alpha_z = 0, H = (H_1, \dots, H_r) \in (\mathbb{R}^{d \times d})^r \right\},$$

equipped with the norm

$$\|(\alpha, H)\|_{\text{tan}} := \|\alpha\|_1 + \|H\|_{\text{lib}}.$$

**Lemma F.6** (Derivative of the characteristic map). Fix  $t \in \mathbb{R}^p$ . The map  $\eta \mapsto \psi_t(\eta)$  is  $C^\infty$  on  $\Xi(B, \mathcal{L})$ . Its Fréchet derivative at  $\eta^* = (\pi^*, \Theta^*)$  in direction  $u = (\alpha, H) \in \mathcal{T}_{\eta^*}$  is

$$D\psi_{t, \eta^*}[u] = \sum_{z \in \mathcal{L}} \alpha_z e^{i\langle t, v_z^* \rangle} + i \sum_{z \in \mathcal{L}} \pi_z^* e^{i\langle t, v_z^* \rangle} \langle t, Dv_{z, \Theta^*}[H] \rangle, \quad (89)$$

where

$$v_z^* := v_z(\Theta^*), \quad Dv_{z, \Theta^*}[H] := \text{vec}([D\Phi_{\Theta^*}(H)]_z).$$

*Proof.* Smoothness follows because  $\eta \mapsto v_z(\Theta)$  is polynomial for each  $z \in \mathcal{L}$ , and  $\eta \mapsto \pi_z e^{i\langle t, v_z(\Theta) \rangle}$  is therefore smooth. Differentiating Equation (83) gives

$$D\psi_{t, \eta^*}[u] = \sum_{z \in \mathcal{L}} \alpha_z e^{i\langle t, v_z^* \rangle} + \sum_{z \in \mathcal{L}} \pi_z^* i e^{i\langle t, v_z^* \rangle} \langle t, Dv_{z, \Theta^*}[H] \rangle,$$

which is Equation (89). □

**Proposition F.7** (Injectivity of the linearized characteristic map). Assume that

$$\gamma_{\text{comp}}(\eta^*) > 0 \quad \text{and} \quad \lambda_{\text{lin}}(\eta^*) > 0. \quad (90)$$

Then the linear map

$$D\Psi_{\eta^*} : \mathcal{T}_{\eta^*} \rightarrow C(\mathbb{R}^p; \mathbb{C}), \quad u \mapsto (t \mapsto D\psi_{t, \eta^*}[u])$$

is injective.

*Proof.* Fix  $u = (\alpha, H) \in \mathcal{T}_{\eta^*}$ , and suppose that

$$D\psi_{t, \eta^*}[u] = 0 \quad \text{for all } t \in \mathbb{R}^p.$$

Define the compactly supported distribution

$$\nu_u := \sum_{z \in \mathcal{L}} \alpha_z \delta_{v_z^*} - \sum_{z \in \mathcal{L}} \pi_z^* \sum_{\ell=1}^p [Dv_{z, \Theta^*}[H]]_\ell \partial_\ell \delta_{v_z^*}. \quad (91)$$

With the Fourier-transform convention

$$\widehat{\delta}_x(t) = e^{i\langle t, x \rangle}, \quad \widehat{\partial_\ell \delta_x}(t) = it_\ell e^{i\langle t, x \rangle},$$

the identity Equation (89) shows that

$$\widehat{\nu}_u(t) = D\psi_{t, \eta^*}[u].$$

Hence  $\widehat{\nu}_u \equiv 0$ , so  $\nu_u = 0$  by injectivity of the Fourier transform on compactly supported distributions.

Because  $\gamma_{\text{comp}}(\eta^*) > 0$ , the support points  $\{v_z^* : z \in \mathcal{L}\}$  are pairwise distinct. Fix  $z_0 \in \mathcal{L}$ , and choose  $\chi_{z_0} \in C_c^\infty(\mathbb{R}^p)$  such that

$$\chi_{z_0}(v_{z_0}^*) = 1, \quad \nabla \chi_{z_0}(v_{z_0}^*) = 0, \quad \chi_{z_0} \equiv 0 \quad \text{in a neighborhood of } v_z^* \text{ for all } z \neq z_0.$$

Pairing  $\nu_u$  with  $\chi_{z_0}$  gives

$$0 = \langle \nu_u, \chi_{z_0} \rangle = \alpha_{z_0}.$$

Since  $z_0$  was arbitrary,  $\alpha_z = 0$  for all  $z \in \mathcal{L}$ .

Next, for each  $\ell \in \{1, \dots, p\}$ , define

$$\chi_{z_0, \ell}(x) := (x_\ell - v_{z_0, \ell}^*) \chi_{z_0}(x).$$

Then

$$\chi_{z_0, \ell}(v_{z_0}^*) = 0, \quad \nabla \chi_{z_0, \ell}(v_{z_0}^*) = e_\ell.$$

Since  $\nu_u = 0$  and  $\alpha = 0$ ,

$$0 = \langle \nu_u, \chi_{z_0, \ell} \rangle = -\pi_{z_0}^* [Dv_{z_0, \Theta^*}[H]]_\ell.$$

Because  $\pi_{z_0}^* > 0$ , we obtain

$$Dv_{z_0, \Theta^*}[H] = 0.$$

As  $z_0$  was arbitrary,

$$D\Phi_{\Theta^*}(H) = 0.$$

Finally,  $\lambda_{\text{lin}}(\eta^*) > 0$  implies

$$0 = \|D\Phi_{\Theta^*}(H)\|_{\pi^*}^2 \geq \lambda_{\text{lin}}(\eta^*) \|H\|_{\text{lib}}^2,$$

hence  $H = 0$ . Therefore  $u = (\alpha, H) = 0$ , proving injectivity.  $\square$

The next proposition converts the infinite-dimensional injective witness  $t \mapsto \psi_t(\eta)$  into a finite-dimensional one.

**Proposition F.8** (Finite witness selection). *Under Equation (90), there exist finitely many frequencies*

$$T_{\text{reg}} = \{t_1, \dots, t_J\} \subset \mathbb{R}^p$$

and a constant  $c_{\text{freq}} > 0$  such that

$$\|DT_{\text{reg}, \eta^*}[u]\|_2 \geq c_{\text{freq}} \|u\|_{\text{tan}} \quad \text{for every } u \in \mathcal{T}_{\eta^*}. \quad (92)$$

*Proof.* Consider the unit sphere

$$\mathbb{S}_{\text{tan}} := \{u \in \mathcal{T}_{\eta^*} : \|u\|_{\text{tan}} = 1\},$$

which is compact because  $\mathcal{T}_{\eta^*}$  is finite-dimensional. By Theorem F.7, for each  $u \in \mathbb{S}_{\text{tan}}$  there exists  $t(u) \in \mathbb{R}^p$  such that

$$|D\psi_{t(u), \eta^*}[u]| > 0.$$

By continuity of  $v \mapsto D\psi_{t(u), \eta^*}[v]$ , there exists an open neighborhood  $U_u$  of  $u$  in  $\mathbb{S}_{\text{tan}}$  such that

$$|D\psi_{t(u), \eta^*}[v]| \geq \frac{1}{2} |D\psi_{t(u), \eta^*}[u]| \quad \text{for all } v \in U_u.$$

The family  $\{U_u : u \in \mathbb{S}_{\text{tan}}\}$  is an open cover of  $\mathbb{S}_{\text{tan}}$ , so compactness yields a finite subcover  $U_{u_1}, \dots, U_{u_J}$ . Set

$$t_j := t(u_j), \quad T_{\text{reg}} := \{t_1, \dots, t_J\}.$$

Define

$$c_{\text{freq}} := \min_{1 \leq j \leq J} \frac{1}{2} |D\psi_{t_j, \eta^*}[u_j]| > 0.$$

Then for every  $u \in \mathbb{S}_{\text{tan}}$ , at least one  $j$  satisfies  $u \in U_{u_j}$ , hence

$$|D\psi_{t_j, \eta^*}[u]| \geq c_{\text{freq}}.$$

Therefore

$$\|DT_{\text{reg}, \eta^*}[u]\|_2 \geq \max_{1 \leq j \leq J} |D\psi_{t_j, \eta^*}[u]| \geq c_{\text{freq}} \quad \text{for all } u \in \mathbb{S}_{\text{tan}}.$$

By homogeneity, Equation (92) follows for all  $u \in \mathcal{T}_{\eta^*}$ .  $\square$

**Theorem F.9** (Regular local inverse inequality). *Assume Equation (90), and let  $T_{\text{reg}}$  be the frequency set from Theorem F.8. Fix  $0 < \delta < 1$ . Then there exist constants*

$$\rho_{\text{reg}} > 0, \quad c_{\text{reg}} > 0, \quad C_{\text{reg}}(\delta) > 0$$

such that on  $\mathcal{G}_{m,n}(\delta)$ , for every  $\eta \in \Xi(B, \mathcal{L})$  with

$$d_q(\eta, \eta^*) \leq \rho_{\text{reg}},$$

one has

$$c_{\text{reg}} d_q(\eta, \eta^*) \leq \|\mathbb{T}_{T_{\text{reg}}}(\eta) - \mathbb{T}_{T_{\text{reg}}}(\eta^*)\|_2 \leq C_{\text{reg}}(\delta) \bar{h}_X(\eta, \eta^*), \quad (93)$$

and consequently

$$\bar{h}_X^2(\eta, \eta^*) \geq c_{\text{inv}}(\delta) d_q(\eta, \eta^*)^2 \quad \text{for all } d_q(\eta, \eta^*) \leq \rho_{\text{reg}}, \quad (94)$$

where  $c_{\text{inv}}(\delta) := c_{\text{reg}}^2 / C_{\text{reg}}(\delta)^2$ .

*Proof.* The upper bound in Equation (93) is exactly Equation (87) applied to  $T = T_{\text{reg}}$ .

For the lower bound, write  $T := T_{\text{reg}}$ . Since each coordinate of  $T$  is  $C^\infty$ , and  $\eta^*$  lies in the interior of the simplex, there exists a neighborhood  $\mathcal{N}^* \subset \Xi(B, \mathcal{L})$  of  $\eta^*$  and a constant  $C_{\text{rem}} < \infty$  such that

$$\|T(\eta^* + u) - T(\eta^*) - DT_{\eta^*}[u]\|_2 \leq C_{\text{rem}} \|u\|_{\text{tan}}^2 \quad (95)$$

whenever  $\eta^* + u \in \mathcal{N}^*$ .

Fix  $\eta$  with  $d_q(\eta, \eta^*)$  sufficiently small. By definition of the quotient metric, there exists  $\sigma \in \text{Aut}(\mathcal{L})$  such that

$$d_q(\eta, \eta^*) = d_{\text{par}}(\sigma \cdot \eta, \eta^*).$$

Set

$$\bar{\eta} := \sigma \cdot \eta, \quad u := \bar{\eta} - \eta^* \in \mathcal{T}_{\eta^*}.$$

Then

$$\|u\|_{\text{tan}} = d_{\text{par}}(\bar{\eta}, \eta^*) = d_q(\eta, \eta^*),$$

and by quotient invariance of  $T$ ,

$$T(\eta) = T(\bar{\eta}).$$

Using Theorem F.8 and Equation (95),

$$\|T(\eta) - T(\eta^*)\|_2 = \|T(\bar{\eta}) - T(\eta^*)\|_2 \geq \|DT_{\eta^*}[u]\|_2 - C_{\text{rem}} \|u\|_{\text{tan}}^2 \geq c_{\text{freq}} \|u\|_{\text{tan}} - C_{\text{rem}} \|u\|_{\text{tan}}^2.$$

Choose

$$\rho_{\text{reg}} := \min \left\{ \text{dist}(\eta^*, \Xi(B, \mathcal{L}) \setminus \mathcal{N}^*), \frac{c_{\text{freq}}}{2C_{\text{rem}}} \right\}, \quad c_{\text{reg}} := \frac{c_{\text{freq}}}{2}.$$

Then for  $d_q(\eta, \eta^*) = \|u\|_{\text{tan}} \leq \rho_{\text{reg}}$ ,

$$\|T(\eta) - T(\eta^*)\|_2 \geq \frac{c_{\text{freq}}}{2} \|u\|_{\text{tan}} = c_{\text{reg}} d_q(\eta, \eta^*).$$

This proves the lower bound in Equation (93). Combining it with the upper bound gives Equation (94).  $\square$

*Remark F.10* (Regular order-1 geometry). Theorem F.9 is the stochastic-task-grammar analogue of strong identifiability in finite mixtures: the predictive geometry is locally Euclidean in the quotient parameter, hence the inverse exponent is 1. This is precisely the regular regime in the sense of singular learning theory and mixture geometry (Watanabe, 2009; Ho & Nguyen, 2016; 2019).

### E.3. Structural posterior contraction under a local inverse inequality

We now isolate the exact abstract implication needed later: predictive contraction plus a local inverse inequality yields structural contraction.

**Theorem F.11** (Structural contraction from predictive contraction). *Fix  $0 < \delta < 1$ , and suppose that Equation (75) holds. Assume that there exist constants*

$$\kappa \geq 1, \quad \rho_\star > 0, \quad c_\star > 0, \quad \underline{h}_\star > 0$$

such that on  $\mathcal{G}_{m,n}(\delta)$ , for all sufficiently large  $m$  and every  $\eta \in \Xi(B, \mathcal{L})$ ,

$$d_q(\eta, \eta^\star) \leq \rho_\star \implies \bar{h}_X^2(\eta, \eta^\star) \geq c_\star d_q(\eta, \eta^\star)^{2\kappa}, \quad (96)$$

$$d_q(\eta, \eta^\star) \geq \rho_\star \implies \bar{h}_X(\eta, \eta^\star) \geq \underline{h}_\star. \quad (97)$$

Then there exists  $M_{\text{str}} \in (0, \infty)$  such that for every  $M \geq M_{\text{str}}$ ,

$$\Pi_m \left( \eta : d_q(\eta, \eta^\star) > M \delta_{m,n}^{1/\kappa} \right) \longrightarrow 0 \quad (98)$$

in  $P_{\eta^\star}$ -probability.

*Proof.* Fix  $M > 0$ , and let

$$A_m(M) := \left\{ \eta : d_q(\eta, \eta^\star) > M \delta_{m,n}^{1/\kappa} \right\}.$$

Write

$$A_m(M) = A_m^{\text{loc}}(M) \cup A_m^{\text{far}}(M),$$

where

$$A_m^{\text{loc}}(M) := A_m(M) \cap \{d_q(\eta, \eta^\star) \leq \rho_\star\}, \quad A_m^{\text{far}}(M) := A_m(M) \cap \{d_q(\eta, \eta^\star) > \rho_\star\}.$$

Because  $\delta_{m,n} \rightarrow 0$ , for every fixed  $M$  we have  $M \delta_{m,n}^{1/\kappa} < \rho_\star$  eventually, so  $A_m^{\text{loc}}(M)$  is nonempty only inside the local inverse regime. On  $\mathcal{G}_{m,n}(\delta)$ , every  $\eta \in A_m^{\text{loc}}(M)$  satisfies by Equation (96)

$$\bar{h}_X^2(\eta, \eta^\star) \geq c_\star M^{2\kappa} \delta_{m,n}^2.$$

Hence

$$A_m^{\text{loc}}(M) \subseteq \left\{ \eta : \bar{h}_X(\eta, \eta^\star) \geq \sqrt{c_\star} M^\kappa \delta_{m,n} \right\}.$$

Likewise, by Equation (97),

$$A_m^{\text{far}}(M) \subseteq \left\{ \eta : \bar{h}_X(\eta, \eta^\star) \geq \underline{h}_\star \right\}.$$

Therefore on  $\mathcal{G}_{m,n}(\delta)$ ,

$$\Pi_m(A_m(M)) \leq \Pi_m(\bar{h}_X(\eta, \eta^\star) \geq \sqrt{c_\star} M^\kappa \delta_{m,n}) + \Pi_m(\bar{h}_X(\eta, \eta^\star) \geq \underline{h}_\star).$$

Choose  $M_{\text{str}}$  so large that

$$\sqrt{c_\star} M^\kappa \geq M_{\text{pred}} \quad \text{for all } M \geq M_{\text{str}},$$

where  $M_{\text{pred}}$  is the predictive-contraction constant from Theorem E.1. Then the first posterior term above converges to 0 by Theorem E.1, while the second also converges to 0 because  $\underline{h}_\star$  is a fixed positive number and  $\delta_{m,n} \rightarrow 0$ . Since  $P_{\eta^\star}(\mathcal{G}_{m,n}(\delta)^c) \rightarrow 0$ , the same conclusion holds unconditionally in  $P_{\eta^\star}$ -probability.  $\square$

**Corollary F.12** (Regular structural contraction). *Assume Equation (90), the design regime Equation (75), and the empirical separation condition*

$$\inf_{\eta: d_q(\eta, \eta^\star) \geq \rho_{\text{reg}}} \bar{h}_X(\eta, \eta^\star) \geq \underline{h}_{\text{reg}} \quad \text{on } \mathcal{G}_{m,n}(\delta) \quad (99)$$

for some  $\underline{h}_{\text{reg}} > 0$  and all sufficiently large  $m$ , where  $\rho_{\text{reg}}$  is the radius from Theorem F.9. Then there exists  $M_{\text{reg}} \in (0, \infty)$  such that for every  $M \geq M_{\text{reg}}$ ,

$$\Pi_m(\eta : d_q(\eta, \eta^\star) > M \delta_{m,n}) \longrightarrow 0 \quad (100)$$

in  $P_{\eta^\star}$ -probability.

*Proof.* Apply Theorem F.11 with  $\kappa = 1$ , local inverse Equation (94), and global separation Equation (99).  $\square$

*Remark F.13* (On the global separation condition). The local inverse inequality is the genuinely nontrivial part of the structural theory. The separation condition Equation (99) is a compactness / identifiability requirement that prevents predictive near-collisions away from the truth. It is automatic on compact quotient classes on which the reduced predictive map is globally identifiable.

#### F.4. An explicit quadratic singular family

We now exhibit a concrete family in which the first-order geometry vanishes and the inverse exponent becomes 2.

**Definition F.14** (Duplicated-module singular family). Let

$$r = 2, \quad \mathcal{L}_{\text{dup}} := \{(1), (2)\}, \quad \pi^* = \left(\frac{1}{2}, \frac{1}{2}\right).$$

Fix  $A^* \in \mathbb{R}^{d \times d}$  with  $\|A^*\|_{\text{op}} < B$ , and define the singular center

$$\eta_0 := (\pi^*, (A^*, A^*)).$$

For any  $\Delta \in \mathbb{R}^{d \times d}$  such that  $\|A^* \pm \Delta\|_{\text{op}} \leq B$ , define

$$\eta_\Delta := (\pi^*, (A^* + \Delta, A^* - \Delta)). \quad (101)$$

**Lemma F.15** (Quotient distance in the duplicated family). *For every  $\Delta$  in the duplicated family,*

$$d_q(\eta_\Delta, \eta_0) = 2\|\Delta\|_F. \quad (102)$$

*Proof.* The automorphism group of  $\mathcal{L}_{\text{dup}}$  consists of the identity and the transposition  $1 \leftrightarrow 2$ . Under either alignment, the weight part contributes 0, while the module part contributes

$$\|A^* + \Delta - A^*\|_F + \|A^* - \Delta - A^*\|_F = 2\|\Delta\|_F.$$

Hence the minimum in the quotient metric equals  $2\|\Delta\|_F$ .  $\square$

Fix a task  $i$ , and write

$$\Sigma_i := (X_i X_i^\top)^{-1} \otimes \sigma^2 I_d, \quad \mu^* := \text{vec}(A^*), \quad s := \text{vec}(\Delta).$$

Conditional on  $X_i$ , the reduced law under  $\eta_0$  is

$$q_{0,i} = \mathcal{N}(\mu^*, \Sigma_i),$$

while under  $\eta_\Delta$  it is the symmetric two-point Gaussian mixture

$$q_{\Delta,i} = \frac{1}{2}\mathcal{N}(\mu^* + s, \Sigma_i) + \frac{1}{2}\mathcal{N}(\mu^* - s, \Sigma_i).$$

**Lemma F.16** (Exact  $\chi^2$ -divergence in the duplicated family). *For every task  $i \in [m]$ ,*

$$r_{i,\Delta}^2 := s^\top \Sigma_i^{-1} s = \frac{1}{\sigma^2} \text{tr}(\Delta X_i X_i^\top \Delta^\top). \quad (103)$$

*Then*

$$\chi^2(q_{\Delta,i}, q_{0,i}) = \cosh(r_{i,\Delta}^2) - 1. \quad (104)$$

*Proof.* Let  $p_0$  and  $p_\Delta$  be the densities of  $q_{0,i}$  and  $q_{\Delta,i}$ . Write

$$u := \Sigma_i^{-1/2}(w - \mu^*), \quad a := \Sigma_i^{-1/2}s, \quad \|a\|_2^2 = r_{i,\Delta}^2.$$

A direct Gaussian calculation gives the likelihood ratio

$$\frac{p_\Delta(w)}{p_0(w)} = \frac{1}{2} \exp\left(a^\top u - \frac{1}{2}\|a\|_2^2\right) + \frac{1}{2} \exp\left(-a^\top u - \frac{1}{2}\|a\|_2^2\right) = e^{-\|a\|_2^2/2} \cosh(a^\top u).$$

Therefore

$$1 + \chi^2(q_{\Delta,i}, q_{0,i}) = \int \frac{p_\Delta(w)^2}{p_0(w)} dw = \mathbb{E}_{U \sim \mathcal{N}(0, I_p)} \left[ e^{-\|a\|_2^2} \cosh^2(a^\top U) \right].$$

Since  $a^\top U \sim \mathcal{N}(0, \|a\|_2^2)$  and

$$\cosh^2 x = \frac{\cosh(2x) + 1}{2},$$

we obtain

$$1 + \chi^2(q_{\Delta,i}, q_{0,i}) = e^{-r_{i,\Delta}^2} \frac{\mathbb{E}[\cosh(2T)] + 1}{2}, \quad T \sim \mathcal{N}(0, r_{i,\Delta}^2).$$

Because  $\mathbb{E}[e^{2T}] = e^{2r_{i,\Delta}^2}$ , we have  $\mathbb{E}[\cosh(2T)] = e^{2r_{i,\Delta}^2}$ . Thus

$$1 + \chi^2(q_{\Delta,i}, q_{0,i}) = \frac{e^{r_{i,\Delta}^2} + e^{-r_{i,\Delta}^2}}{2} = \cosh(r_{i,\Delta}^2),$$

which proves Equation (104). □

**Proposition F.17** (Quadratic flattening in the duplicated family). *Fix  $0 < \delta < 1$ . There exist constants*

$$\rho_{\text{sing}}(\delta) > 0, \quad 0 < c_{\text{sing}}(\delta) \leq C_{\text{sing}}(\delta) < \infty$$

such that on  $\mathcal{G}_{m,n}(\delta)$ , for every  $\Delta$  satisfying

$$\sqrt{n} \|\Delta\|_F \leq \rho_{\text{sing}}(\delta),$$

one has

$$c_{\text{sing}}(\delta) n^2 d_q(\eta_\Delta, \eta_0)^4 \leq \bar{h}_X^2(\eta_\Delta, \eta_0) \leq C_{\text{sing}}(\delta) n^2 d_q(\eta_\Delta, \eta_0)^4. \quad (105)$$

In particular, the duplicated family has local inverse exponent 2.

*Proof.* Fix  $i \in [m]$ . On  $\mathcal{G}_{m,n}(\delta)$ ,

$$(1 - \delta)nI_d \preceq X_i X_i^\top \preceq (1 + \delta)nI_d,$$

hence by Equation (103),

$$\frac{(1 - \delta)n}{\sigma^2} \|\Delta\|_F^2 \leq r_{i,\Delta}^2 \leq \frac{(1 + \delta)n}{\sigma^2} \|\Delta\|_F^2. \quad (106)$$

Choose

$$\rho_{\text{sing}}(\delta) := \frac{\sigma}{\sqrt{1 + \delta}},$$

so that  $\sqrt{n} \|\Delta\|_F \leq \rho_{\text{sing}}(\delta)$  implies  $r_{i,\Delta}^2 \leq 1$  for all  $i$ .

**Upper bound.** For any pair of dominated laws,  $h^2(P, Q) \leq \chi^2(P, Q)/2$ . Thus by Theorem F.16,

$$h^2(q_{\Delta,i}, q_{0,i}) \leq \frac{1}{2} (\cosh(r_{i,\Delta}^2) - 1).$$

Since  $r_{i,\Delta}^2 \leq 1$ , Taylor's theorem gives

$$\cosh(u) - 1 \leq \frac{e}{2} u^2 \quad \text{for } u \in [0, 1].$$

Therefore

$$h^2(q_{\Delta,i}, q_{0,i}) \leq \frac{e}{4} r_{i,\Delta}^4 \leq \frac{e(1 + \delta)^2}{4\sigma^4} n^2 \|\Delta\|_F^4.$$

**Lower bound.** Let

$$a_i := \Sigma_i^{-1/2} s \in \mathbb{R}^p, \quad r_i := \|a_i\|_2 = r_{i,\Delta}.$$

Under  $q_{0,i}$ , the whitened variable

$$U := \Sigma_i^{-1/2} (\text{vec}(\widehat{W}_i) - \mu^*)$$

has law  $\mathcal{N}(0, I_p)$ . Under  $q_{\Delta,i}$ ,  $U$  has the symmetric mixture

$$\frac{1}{2}\mathcal{N}(a_i, I_p) + \frac{1}{2}\mathcal{N}(-a_i, I_p).$$

Define the quadratic witness

$$f_i(u) := (a_i^\top u)^2 - r_i^2.$$

Under  $q_{0,i}$ ,  $a_i^\top U \sim \mathcal{N}(0, r_i^2)$ , so

$$\mathbb{E}_{q_{0,i}}[f_i(U)] = 0, \quad \mathbb{E}_{q_{0,i}}[f_i(U)^2] = 2r_i^4.$$

Under  $q_{\Delta,i}$ , a direct Gaussian-moment computation yields

$$\mathbb{E}_{q_{\Delta,i}}[f_i(U)] = r_i^4, \quad \mathbb{E}_{q_{\Delta,i}}[f_i(U)^2] = 2r_i^4 + 4r_i^6 + r_i^8 \leq 7r_i^4,$$

since  $r_i^2 \leq 1$ .

For any square-integrable  $f$  and any two dominated laws  $P, Q$ ,

$$|E_P f - E_Q f| = \left| \int f(\sqrt{p} - \sqrt{q})(\sqrt{p} + \sqrt{q}) \right| \leq 2\sqrt{E_P f^2 + E_Q f^2} h(P, Q).$$

Applying this with  $P = q_{\Delta,i}$ ,  $Q = q_{0,i}$ , and  $f = f_i$ , we obtain

$$r_i^4 = |\mathbb{E}_{q_{\Delta,i}}[f_i] - \mathbb{E}_{q_{0,i}}[f_i]| \leq 2\sqrt{9r_i^4} h(q_{\Delta,i}, q_{0,i}) = 6r_i^2 h(q_{\Delta,i}, q_{0,i}).$$

Hence

$$h^2(q_{\Delta,i}, q_{0,i}) \geq \frac{r_i^4}{36}.$$

Using the lower side of Equation (106),

$$h^2(q_{\Delta,i}, q_{0,i}) \geq \frac{(1-\delta)^2}{36\sigma^4} n^2 \|\Delta\|_F^4.$$

Averaging over  $i = 1, \dots, m$  and recalling Equation (102),

$$\bar{h}_X^2(\eta_\Delta, \eta_0) = \frac{1}{m} \sum_{i=1}^m h^2(q_{\Delta,i}, q_{0,i}) \asymp n^2 \|\Delta\|_F^4 \asymp n^2 d_q(\eta_\Delta, \eta_0)^4,$$

with constants depending only on  $(\delta, \sigma)$ . This proves Equation (105).  $\square$

**Corollary F.18** (Quadratic singular structural contraction on the duplicated family). *Assume that the prior is supported on the duplicated family  $\{\eta_\Delta : \sqrt{n} \|\Delta\|_F \leq \rho_{\text{sing}}(\delta)\}$ , and that its induced density on  $\Delta$  is continuous and strictly positive near  $\Delta = 0$ . Then there exists  $M_{\text{dup}} \in (0, \infty)$  such that for every  $M \geq M_{\text{dup}}$ ,*

$$\Pi_m \left( \eta_\Delta : d_q(\eta_\Delta, \eta_0) > M \frac{\delta_{m,n}^{1/2}}{\sqrt{n}} \right) \rightarrow 0 \quad (107)$$

in  $P_{\eta_0}$ -probability.

*Proof.* On the support of the prior, Theorem F.17 supplies the local inverse inequality

$$\bar{h}_X^2(\eta_\Delta, \eta_0) \geq c_{\text{sing}}(\delta) n^2 d_q(\eta_\Delta, \eta_0)^4.$$

Hence, if

$$d_q(\eta_\Delta, \eta_0) > M \frac{\delta_{m,n}^{1/2}}{\sqrt{n}},$$

then

$$\bar{h}_X(\eta_\Delta, \eta_0) > \sqrt{c_{\text{sing}}(\delta)} M^2 \delta_{m,n}.$$

Choosing  $M_{\text{dup}}$  so that  $\sqrt{c_{\text{sing}}(\delta)} M^2 \geq M_{\text{pred}}$  and applying Theorem E.1 proves Equation (107).  $\square$

*Remark F.19* (Why the duplicated family matters). The duplicated family is the exact stochastic-task-grammar analogue of an overfitted symmetric finite mixture: the first-order perturbation cancels, the predictive law only changes at second order, and the structural inverse exponent becomes 2. The mechanism is classical in singular mixture geometry (Ho & Nguyen, 2016; 2019; Watanabe, 2009), but the calculation above is specific to the reduced grouped experiment induced by compositional task grammars.

## G. Minimax Lower Bounds

This appendix proves four lower-bound statements. First, we establish a grammar-law lower bound of order  $\sqrt{K_{\text{eff}}/m}$  by reducing to a multinomial subproblem. Second, under a natural anchored-language condition, we prove a module-estimation lower bound of order  $d\sqrt{s/(mn)}$  for  $s$  anchor modules. Third, we show that zero contextual observability induces an exact non-identifiability floor. Fourth, on the duplicated-module singular family of Appendix F, we derive a sharp two-point lower bound of order  $m^{-1/4}n^{-1/2}$ , matching the singular exponent  $\kappa = 2$  up to logarithmic factors.

The lower bounds are stated for the full grouped-data experiment  $(X_i, Y_i)_{i=1}^m$ . Whenever convenient, we pass to a more informative oracle experiment or to the sufficient-statistic reduction of Appendix B; this is legitimate because lower bounds are monotone under Blackwell domination.

### G.1. Decision-theoretic setup and generic comparison lemmas

For a parameter class  $\mathcal{H} \subseteq \Xi(B, \mathcal{L})$ , define the minimax structural risk

$$\mathfrak{R}_{m,n}(\mathcal{H}; d_q) := \inf_{\hat{\eta}} \sup_{\eta \in \mathcal{H}} \mathbb{E}_\eta[d_q(\hat{\eta}, \eta)], \quad (108)$$

where the infimum is taken over all measurable estimators  $\hat{\eta} = \hat{\eta}(X_{1:m}, Y_{1:m})$  with values in  $\Xi(B, \mathcal{L})$ .

**Lemma G.1** (More informative experiments can only decrease minimax risk). *Let  $\mathcal{E}_1 = \{P_\eta : \eta \in \mathcal{H}\}$  and  $\mathcal{E}_2 = \{Q_\eta : \eta \in \mathcal{H}\}$  be two statistical experiments on measurable spaces  $(\mathcal{Y}, \mathcal{Y})$  and  $(\mathcal{X}, \mathcal{X})$ , respectively. Suppose that there exists a Markov kernel  $K(dy | x)$ , independent of  $\eta$ , such that*

$$P_\eta(A) = \int_{\mathcal{X}} K(A | x) Q_\eta(dx) \quad \text{for every } A \in \mathcal{Y}, \eta \in \mathcal{H}.$$

Then for every loss  $L : \mathcal{H} \times \Xi(B, \mathcal{L}) \rightarrow [0, \infty)$ ,

$$\inf_{\hat{\eta}_1} \sup_{\eta \in \mathcal{H}} \mathbb{E}_{P_\eta}[L(\eta, \hat{\eta}_1)] \geq \inf_{\hat{\eta}_2} \sup_{\eta \in \mathcal{H}} \mathbb{E}_{Q_\eta}[L(\eta, \hat{\eta}_2)].$$

*Proof.* Fix any estimator  $\hat{\eta}_1 : \mathcal{Y} \rightarrow \Xi(B, \mathcal{L})$  for the less informative experiment  $\mathcal{E}_1$ . Define the randomized estimator on  $\mathcal{E}_2$

$$\hat{\eta}_2(x) \sim K(dy | x) \text{ followed by } \hat{\eta}_1(y).$$

Then for every  $\eta \in \mathcal{H}$ ,

$$\mathbb{E}_{Q_\eta}[L(\eta, \hat{\eta}_2)] = \mathbb{E}_{P_\eta}[L(\eta, \hat{\eta}_1)].$$

Taking the supremum over  $\eta$ , then the infimum over  $\hat{\eta}_1$ , proves the claim.  $\square$

**Corollary G.2** (Two useful monotonicity reductions). *For every parameter class  $\mathcal{H} \subseteq \Xi(B, \mathcal{L})$ :*

1. *the oracle experiment that augments the full grouped data with the latent programs  $Z_{1:m}$  is more informative than the original grouped-data experiment, hence any minimax lower bound proved in the oracle experiment is valid for the original problem;*
2. *the full grouped-data experiment and the reduced experiment  $(X_i, \widehat{W}_i)_{i=1}^m$  are Blackwell equivalent, because Appendix B shows that*

$$Y_i = \widehat{W}_i X_i + R_i,$$

where, conditional on  $(X_i, \widehat{W}_i)$ , the residual  $R_i$  has a law that does not depend on  $\eta$ .

*Proof.* Item (1) is immediate from Theorem G.1: forgetting the latent labels  $Z_{1:m}$  maps the oracle experiment to the original one.

For item (2), the map  $(X_i, Y_i) \mapsto (X_i, \widehat{W}_i)$  is measurable, so the reduced experiment is less informative than the full one. Conversely, by Theorems B.5 and B.6 in Appendix B, conditional on  $(X_i, \widehat{W}_i)$ , the residual  $R_i$  is independent of  $\eta$  with parameter-free law

$$R_i \mid (X_i, \widehat{W}_i) \sim \mathcal{MN}_{d,n}(0, \sigma^2 I_d, I_n - P_{X_i}).$$

Drawing  $R_i$  from this law and setting  $Y_i = \widehat{W}_i X_i + R_i$  reconstructs the full conditional law of  $Y_i$ . Hence the full and reduced experiments dominate one another.  $\square$

We also use the following classical decision-theoretic tools.

**Lemma G.3** (Le Cam's two-point method). *Let  $\eta_0, \eta_1 \in \mathcal{H}$ , and let  $P_0, P_1$  denote the corresponding laws of the observed experiment. Then*

$$\inf_{\widehat{\eta}} \sup_{\eta \in \{\eta_0, \eta_1\}} \mathbb{E}_{\eta} [d_q(\widehat{\eta}, \eta)] \geq \frac{d_q(\eta_0, \eta_1)}{4} \left(1 - \text{TV}(P_0, P_1)\right). \quad (109)$$

Moreover,

$$\text{TV}(P_0, P_1) \leq \sqrt{\frac{1}{2} \chi^2(P_1, P_0)}.$$

*Proof.* Let

$$r := \frac{1}{2} d_q(\eta_0, \eta_1), \quad A := \{\widehat{\eta} : d_q(\widehat{\eta}, \eta_0) < r\}.$$

By the triangle inequality, on  $A$  we have  $d_q(\widehat{\eta}, \eta_1) \geq r$ , while on  $A^c$  we have  $d_q(\widehat{\eta}, \eta_0) \geq r$ . Therefore

$$\mathbb{E}_{\eta_0} [d_q(\widehat{\eta}, \eta_0)] + \mathbb{E}_{\eta_1} [d_q(\widehat{\eta}, \eta_1)] \geq r (P_0(A^c) + P_1(A)).$$

Taking the maximum of the two expectations yields

$$\sup_{\eta \in \{\eta_0, \eta_1\}} \mathbb{E}_{\eta} [d_q(\widehat{\eta}, \eta)] \geq \frac{r}{2} (P_0(A^c) + P_1(A)).$$

Finally,

$$\inf_A (P_0(A^c) + P_1(A)) = 1 - \text{TV}(P_0, P_1),$$

which gives Equation (109). The  $\chi^2$ -to-TV bound is standard:

$$\text{TV}(P_0, P_1) = \frac{1}{2} \int |p_1 - p_0| \leq \frac{1}{2} \left( \int \frac{(p_1 - p_0)^2}{p_0} \right)^{1/2} = \sqrt{\frac{1}{2} \chi^2(P_1, P_0)}.$$

$\square$

**Lemma G.4** (Fano's inequality; e.g., Tsybakov (2009, Theorem 2.5)). *Let  $\eta^0, \eta^1, \dots, \eta^M \in \mathcal{H}$ , with  $M \geq 2$ , and suppose that*

$$d_q(\eta^j, \eta^k) \geq 2s \quad \text{for all } 0 \leq j < k \leq M.$$

*If  $P_j$  denotes the law of the observed experiment under  $\eta^j$ , then every estimator  $\hat{\eta}$  satisfies*

$$\sup_{0 \leq j \leq M} \mathbb{E}_{\eta^j} [d_q(\hat{\eta}, \eta^j)] \geq s \left( 1 - \frac{\frac{1}{M} \sum_{j=1}^M \text{KL}(P_j, P_0) + \log 2}{\log M} \right). \quad (110)$$

**Remark G.5** (Varshamov–Gilbert packing). We repeatedly use the standard coding-theoretic fact that for every integer  $q \geq 8$ , there exists a subset  $V \subseteq \{-1, 1\}^q$  such that

$$\log |V| \geq \frac{q}{8} \quad \text{and} \quad d_H(u, v) \geq \frac{q}{8} \quad \text{for all distinct } u, v \in V.$$

This is the usual Varshamov–Gilbert bound; see, e.g., Tsybakov (2009, Lemma 2.9).

## G.2. A generic asymmetric baseline library

The quotient metric identifies parameters up to automorphisms of the admissible language. To construct lower-bound submodels, it is convenient to freeze a baseline library whose automorphism stabilizer is trivial.

**Lemma G.6** (Explicit asymmetric baseline library). *Define*

$$A_a^\circ := \frac{Ba}{4r} I_d, \quad a \in [r]. \quad (111)$$

*Then  $A_a^\circ \in \mathbb{A}_B$  for every  $a \in [r]$ , and*

$$\vartheta_0 := \min_{\sigma \in \text{Aut}(\mathcal{L}) \setminus \{\text{id}\}} \sum_{a=1}^r \|A_a^\circ - A_{\sigma^{-1}(a)}^\circ\|_F > 0. \quad (112)$$

*Consequently, for every submodel that keeps the grammar weights fixed and perturbs the library by at most  $\vartheta_0/8$  in  $d_{\text{par}}$ , the minimizing automorphism in the quotient metric is the identity.*

*Proof.* Since

$$\|A_a^\circ\|_{\text{op}} = \frac{Ba}{4r} \leq \frac{B}{4} < B,$$

we have  $A_a^\circ \in \mathbb{A}_B$ . If  $\sigma \neq \text{id}$ , then there exists at least one  $a \in [r]$  with  $\sigma^{-1}(a) \neq a$ . Therefore

$$\|A_a^\circ - A_{\sigma^{-1}(a)}^\circ\|_F = \frac{B|a - \sigma^{-1}(a)|}{4r} \|I_d\|_F = \frac{B\sqrt{d}}{4r} |a - \sigma^{-1}(a)| > 0.$$

Summing over  $a$  and minimizing over the finite set  $\text{Aut}(\mathcal{L}) \setminus \{\text{id}\}$  yields  $\vartheta_0 > 0$ .

Now let  $\eta, \eta'$  be two parameters in a fixed-weight submodel with

$$d_{\text{par}}(\eta, \eta^\circ) \leq \frac{\vartheta_0}{8}, \quad d_{\text{par}}(\eta', \eta^\circ) \leq \frac{\vartheta_0}{8},$$

where  $\eta^\circ$  has baseline library  $\Theta^\circ = (A_1^\circ, \dots, A_r^\circ)$ . If  $\sigma \neq \text{id}$ , then by the triangle inequality,

$$d_{\text{par}}(\eta, \sigma \cdot \eta') \geq \sum_{a=1}^r \|A_a^\circ - A_{\sigma^{-1}(a)}^\circ\|_F - d_{\text{par}}(\eta, \eta^\circ) - d_{\text{par}}(\eta', \eta^\circ) \geq \vartheta_0 - \frac{\vartheta_0}{4} = \frac{3\vartheta_0}{4}.$$

On the other hand, the identity alignment yields distance at most  $\vartheta_0/4$ . Hence the quotient minimum is attained by the identity.  $\square$

### G.3. Regular grammar-law lower bound

Let

$$N := |\mathcal{L}|, \quad K_{\text{eff}} = N - 1.$$

Assume  $N \geq 2$ , since the case  $K_{\text{eff}} = 0$  is trivial.

Enumerate the admissible language as

$$\mathcal{L} = \{z_1, \dots, z_N\},$$

and let

$$s := \left\lfloor \frac{N}{2} \right\rfloor.$$

Fix the asymmetric baseline library  $\Theta^\circ$  from Theorem G.6. For  $\alpha \in (0, 1/4]$  and  $u = (u_1, \dots, u_s) \in \{-1, 1\}^s$ , define the grammar law

$$\pi_{z_{2j-1}}^u := \frac{1 + \alpha u_j}{N}, \quad \pi_{z_{2j}}^u := \frac{1 - \alpha u_j}{N}, \quad j = 1, \dots, s, \quad (113)$$

and, if  $N$  is odd, set  $\pi_{z_N}^u = 1/N$ . Let

$$\eta^u := (\pi^u, \Theta^\circ).$$

**Theorem G.7** (Grammar-law minimax lower bound). *There exists a constant  $c_{\text{gram}} > 0$ , depending only on  $(\mathcal{L}, r, d, B)$ , such that for all sufficiently large  $m$ ,*

$$\mathfrak{R}_{m,n}(\{\eta^u : u \in \{-1, 1\}^s\}; d_q) \geq c_{\text{gram}} \sqrt{\frac{K_{\text{eff}}}{m}}. \quad (114)$$

Consequently, the same lower bound holds for any larger parameter class containing this submodel.

*Proof.* By Theorem G.6, for all sufficiently large  $m$  we may choose

$$\alpha_m := c_0 \sqrt{\frac{N}{m}}$$

with  $c_0 > 0$  small enough so that  $\alpha_m \leq 1/4$  and every pairwise  $\ell_1$ -distance  $\|\pi^u - \pi^v\|_1$  is strictly smaller than  $\vartheta_0/4$ . Hence, on the submodel  $\{\eta^u\}$ , the quotient metric is simply

$$d_q(\eta^u, \eta^v) = \|\pi^u - \pi^v\|_1.$$

Choose a Varshamov–Gilbert subset  $V \subseteq \{-1, 1\}^s$  satisfying

$$\log |V| \geq \frac{s}{8}, \quad d_H(u, v) \geq \frac{s}{8} \text{ for all distinct } u, v \in V.$$

For  $u, v \in V$ ,

$$d_q(\eta^u, \eta^v) = \|\pi^u - \pi^v\|_1 = \frac{4\alpha_m}{N} d_H(u, v) \geq \frac{\alpha_m s}{2N}.$$

Since  $s/N \geq 1/3$  for all  $N \geq 2$ ,

$$d_q(\eta^u, \eta^v) \geq \frac{\alpha_m}{6} = \frac{c_0}{6} \sqrt{\frac{N}{m}} \asymp \sqrt{\frac{K_{\text{eff}}}{m}}. \quad (115)$$

We now compare the original grouped-data experiment with the oracle experiment that reveals the latent programs  $Z_{1:m}$ . By Theorem G.2, this oracle experiment is more informative, so it suffices to prove the lower bound there. In the oracle experiment,

$$Z_1, \dots, Z_m \stackrel{\text{i.i.d.}}{\sim} \pi^u.$$

Thus the law is simply the multinomial product measure  $(\pi^u)^{\otimes m}$ .

Fix  $u, v \in V$ . If  $h = d_H(u, v)$ , then the single-sample KL divergence is

$$\begin{aligned} \text{KL}(\pi^u, \pi^v) &= \sum_{j: u_j \neq v_j} \left[ \frac{1 + \alpha_m}{N} \log \frac{1 + \alpha_m}{1 - \alpha_m} + \frac{1 - \alpha_m}{N} \log \frac{1 - \alpha_m}{1 + \alpha_m} \right] \\ &= \frac{2\alpha_m h}{N} \log \frac{1 + \alpha_m}{1 - \alpha_m}. \end{aligned} \quad (116)$$

Since  $\alpha_m \leq 1/4$ , the inequality

$$\log \frac{1 + \alpha_m}{1 - \alpha_m} \leq 4\alpha_m$$

gives

$$\text{KL}(\pi^u, \pi^v) \leq \frac{8\alpha_m^2 h}{N} \leq 8\alpha_m^2 \frac{s}{N} \leq 4\alpha_m^2.$$

Therefore the  $m$ -sample KL divergence in the oracle experiment satisfies

$$\text{KL}((\pi^u)^{\otimes m}, (\pi^v)^{\otimes m}) \leq 4m\alpha_m^2 = 4c_0^2 N.$$

Since  $\log |V| \geq s/8$  and  $s \asymp N$ , choosing  $c_0 > 0$  sufficiently small ensures that

$$\max_{u \in V \setminus \{u^0\}} \text{KL}((\pi^u)^{\otimes m}, (\pi^{u^0})^{\otimes m}) \leq \frac{1}{16} \log |V|$$

for any fixed  $u^0 \in V$ . Applying Theorem G.4 and the separation bound Equation (115) yields

$$\inf_{\hat{\eta}} \sup_{u \in V} \mathbb{E}_u [d_q(\hat{\eta}, \eta^u)] \geq c \sqrt{\frac{N}{m}} \asymp \sqrt{\frac{K_{\text{eff}}}{m}}$$

for a constant  $c > 0$  depending only on  $\mathcal{L}$ . By Theorem G.1, the same lower bound holds for the original grouped-data experiment.  $\square$

*Remark G.8 (Interpretation).* Theorem G.7 shows that even if the module library is known perfectly, the latent task-family law  $\pi$  cannot be estimated faster than  $\sqrt{K_{\text{eff}}/m}$ . This is the irreducible grammar-complexity term in the problem.

#### G.4. Anchored-module lower bound

We now show that when the benchmark explicitly exposes some modules in isolation, one also recovers the  $1/\sqrt{mn}$ -type module-estimation difficulty.

**Definition G.9** (Anchor-language condition). Let  $a_1, \dots, a_s \in [r]$  be distinct symbols. We say that the admissible language  $\mathcal{L}$  has an  $s$ -anchor set if:

1. the singleton programs  $(a_1), \dots, (a_s)$  all belong to  $\mathcal{L}$ ;
2. no program in  $\mathcal{L} \setminus \{(a_1), \dots, (a_s)\}$  contains any of the symbols  $a_1, \dots, a_s$ .

Assume Theorem G.9 with anchor set  $a_1, \dots, a_s$ . Fix a grammar law  $\pi^\circ \in \mathfrak{P}_{\mathcal{L}}^\circ$  and set

$$\beta_j := \pi_{(a_j)}^\circ > 0, \quad \beta_{\max} := \max_{1 \leq j \leq s} \beta_j.$$

Fix the asymmetric baseline library  $\Theta^\circ$  from Theorem G.6, and let

$$p := d^2.$$

Choose an orthonormal basis  $E_1, \dots, E_p$  of  $\mathbb{R}^{d \times d}$  under the Frobenius inner product.

For  $u = (u_{j\ell}) \in \{-1, 1\}^{sp}$  and  $\tau > 0$ , define

$$A_{a_j}^{(u)} := A_{a_j}^\circ + \tau \sum_{\ell=1}^p u_{j\ell} E_\ell, \quad j = 1, \dots, s, \quad (117)$$

and keep all non-anchor modules fixed at their baseline values. Let  $\eta^u$  denote the resulting parameter with grammar law  $\pi^\circ$ .

**Theorem G.10** (Anchored-module minimax lower bound). *Assume that  $\mathcal{L}$  has an  $s$ -anchor set. Then there exists a constant  $c_{\text{anc}} > 0$ , depending only on  $(\mathcal{L}, r, d, B, \sigma, \pi^\circ)$ , such that for all sufficiently large  $m$ ,*

$$\mathfrak{R}_{m,n}(\{\eta^u : u \in \{-1, 1\}^{sp}\}; d_q) \geq c_{\text{anc}} \sigma d \sqrt{\frac{s}{mn}}. \quad (118)$$

*Proof.* Choose

$$\tau_{m,n} := c_1 \frac{\sigma}{\sqrt{mn}}$$

with  $c_1 > 0$  sufficiently small so that (i) all perturbed anchor modules remain in  $\mathbb{A}_B$ , and (ii) the total library perturbation is at most  $\vartheta_0/8$ , where  $\vartheta_0$  is the asymmetry gap from Equation (112). Then by Theorem G.6, the quotient minimum on this submodel is attained by the identity, so for any  $u, v$ ,

$$d_q(\eta^u, \eta^v) = \sum_{j=1}^s \|A_{a_j}^{(u)} - A_{a_j}^{(v)}\|_F.$$

If  $h = d_H(u, v)$  denotes the Hamming distance in  $\{-1, 1\}^{sp}$ , then

$$\sum_{j=1}^s \|A_{a_j}^{(u)} - A_{a_j}^{(v)}\|_F \geq \left( \sum_{j=1}^s \|A_{a_j}^{(u)} - A_{a_j}^{(v)}\|_F^2 \right)^{1/2} = 2\tau_{m,n} \sqrt{h}.$$

Choose a Varshamov–Gilbert subset  $V \subseteq \{-1, 1\}^{sp}$  with

$$\log |V| \geq \frac{sp}{8}, \quad d_H(u, v) \geq \frac{sp}{8} \text{ for all } u \neq v \in V.$$

Hence for distinct  $u, v \in V$ ,

$$d_q(\eta^u, \eta^v) \geq 2\tau_{m,n} \sqrt{\frac{sp}{8}} = c \tau_{m,n} d \sqrt{s} \asymp \sigma d \sqrt{\frac{s}{mn}}. \quad (119)$$

We now compare with the oracle experiment that reveals the latent programs  $Z_{1:m}$ ; by Theorem G.2, it suffices to work there. In the oracle experiment, the augmented one-task law is

$$(Z_i, X_i, Y_i), \quad Z_i \sim \pi^\circ.$$

By the anchor-language condition, perturbing an anchor module  $A_{a_j}$  affects only tasks with  $Z_i = (a_j)$ , because no other program contains  $a_j$ . Therefore, for  $u, v \in V$ , the one-task KL divergence in the oracle experiment equals

$$\begin{aligned} \text{KL}(P_u^{\text{or}}, P_v^{\text{or}}) &= \sum_{j=1}^s \beta_j \cdot \frac{n}{2\sigma^2} \|A_{a_j}^{(u)} - A_{a_j}^{(v)}\|_F^2 \\ &\leq \frac{n\beta_{\max}}{2\sigma^2} \sum_{j=1}^s \|A_{a_j}^{(u)} - A_{a_j}^{(v)}\|_F^2 \\ &= \frac{2n\beta_{\max}\tau_{m,n}^2}{\sigma^2} d_H(u, v). \end{aligned} \quad (120)$$

Thus the  $m$ -task oracle KL divergence is bounded by

$$\text{KL}((P_u^{\text{or}})^{\otimes m}, (P_v^{\text{or}})^{\otimes m}) \leq \frac{2mn\beta_{\max}\tau_{m,n}^2}{\sigma^2} d_H(u, v) \leq 2c_1^2 \beta_{\max} sp.$$

Choosing  $c_1$  sufficiently small ensures that this is at most  $\frac{1}{16} \log |V|$  for all  $u, v \in V$ . Fano's lemma Theorem G.4, together with the separation bound Equation (119), yields

$$\inf_{\hat{\eta}} \sup_{u \in V} \mathbb{E}_u[d_q(\hat{\eta}, \eta^u)] \geq c' \sigma d \sqrt{\frac{s}{mn}}$$

for some  $c' > 0$ . By Theorem G.1, the same lower bound holds in the original grouped-data experiment.  $\square$

*Remark G.11* (Benchmark-design meaning of anchors). Theorem G.10 isolates the statistical value of benchmark cells that expose modules in isolation. Without such anchor tasks, module-space lower bounds depend on more complicated local inverse geometry; with anchors, one recovers the classical  $1/\sqrt{mn}$  parametric difficulty of estimating a shared linear operator.

### G.5. Zero-observability impossibility

The next statement turns exact observational collapse into a nonvanishing minimax floor.

**Theorem G.12** (Exact non-identifiability induces a positive minimax floor). *Suppose there exist two parameters  $\eta_0, \eta_1 \in \Xi(B, \mathcal{L})$  such that*

$$d_q(\eta_0, \eta_1) = \Delta_0 > 0 \quad \text{and} \quad Q_{\eta_0}^{(n)} = Q_{\eta_1}^{(n)} \text{ for every } n \geq 1.$$

*Then for every  $m \geq 1$  and every  $n \geq d + 1$ ,*

$$\mathfrak{R}_{m,n}(\{\eta_0, \eta_1\}; d_q) \geq \frac{\Delta_0}{4}. \quad (121)$$

*In particular, whenever Appendix C produces such a pair through zero contextual observability in a single-occurrence language, the structural risk cannot vanish.*

*Proof.* If the one-task grouped laws are equal, then the  $m$ -task laws are also equal:

$$(Q_{\eta_0}^{(n)})^{\otimes m} = (Q_{\eta_1}^{(n)})^{\otimes m}.$$

Therefore the total variation distance between the two  $m$ -task laws is zero. Applying Le Cam's inequality Equation (109) gives

$$\inf_{\hat{\eta}} \sup_{\eta \in \{\eta_0, \eta_1\}} \mathbb{E}_{\eta} [d_q(\hat{\eta}, \eta)] \geq \frac{d_q(\eta_0, \eta_1)}{4} = \frac{\Delta_0}{4}.$$

□

### G.6. Singular lower bound on the duplicated-module family

We finally prove the lower bound corresponding to the quadratic singular geometry of Appendix F. Recall the duplicated family

$$\eta_t := \eta_{\Delta_t}, \quad \Delta_t := t e_1 e_1^\top, \quad t \geq 0,$$

built around the singular center  $\eta_0$  in Theorem F.14. By Theorem F.15,

$$d_q(\eta_t, \eta_0) = 2t.$$

Let  $\tilde{Q}_t^{(n)}$  denote the reduced one-task predictive law from Appendix E under  $\eta_t$ .

**Lemma G.13** (Exact unconditional  $\chi^2$  divergence in the rank-one duplicated subfamily). *Let*

$$a_t := \frac{t^2}{\sigma^2}.$$

*If  $a_t < 1/2$ , then*

$$\chi^2(\tilde{Q}_t^{(n)}, \tilde{Q}_0^{(n)}) = \frac{1}{2}(1 - 2a_t)^{-n/2} + \frac{1}{2}(1 + 2a_t)^{-n/2} - 1. \quad (122)$$

*Moreover, there exists a universal constant  $C_\chi \in (0, \infty)$  such that whenever*

$$a_t \leq \frac{1}{4n},$$

*one has*

$$\chi^2(\tilde{Q}_t^{(n)}, \tilde{Q}_0^{(n)}) \leq C_\chi n^2 a_t^2 = C_\chi \frac{n^2 t^4}{\sigma^4}. \quad (123)$$

*Proof.* By Appendix F, conditional on  $X$ , the reduced one-task  $\chi^2$ -divergence satisfies

$$\chi^2(q_{t,X}, q_{0,X}) = \cosh(r_{X,t}^2) - 1, \quad r_{X,t}^2 = \frac{1}{\sigma^2} \|\Delta_t X\|_F^2.$$

Since  $\Delta_t = t e_1 e_1^\top$ ,

$$\|\Delta_t X\|_F^2 = t^2 \sum_{j=1}^n X_{1j}^2,$$

and  $\sum_{j=1}^n X_{1j}^2 \sim \chi_n^2$ . Because the design law is the same under both parameters,

$$\chi^2(\tilde{Q}_t^{(n)}, \tilde{Q}_0^{(n)}) = \mathbb{E}_X [\chi^2(q_{t,X}, q_{0,X})] = \mathbb{E} [\cosh(a_t \chi_n^2) - 1].$$

Using

$$\cosh(u) = \frac{e^u + e^{-u}}{2}$$

and the mgf of  $\chi_n^2$ ,

$$\mathbb{E}[e^{\lambda \chi_n^2}] = (1 - 2\lambda)^{-n/2} \quad \text{for } \lambda < 1/2,$$

we obtain

$$\chi^2(\tilde{Q}_t^{(n)}, \tilde{Q}_0^{(n)}) = \frac{1}{2}(1 - 2a_t)^{-n/2} + \frac{1}{2}(1 + 2a_t)^{-n/2} - 1,$$

which is Equation (122).

For the upper bound, define

$$f(a) := \frac{1}{2}(1 - 2a)^{-n/2} + \frac{1}{2}(1 + 2a)^{-n/2} - 1.$$

Then  $f(0) = f'(0) = 0$ , and

$$f''(a) = \frac{n(n+2)}{2} \left[ (1 - 2a)^{-n/2-2} + (1 + 2a)^{-n/2-2} \right].$$

If  $0 \leq a \leq 1/(4n)$ , then

$$1 - 2a \geq 1 - \frac{1}{2n}, \quad 1 + 2a \leq 1 + \frac{1}{2n},$$

so both factors above are bounded by a universal constant. Therefore

$$\sup_{0 \leq a \leq 1/(4n)} |f''(a)| \leq C_\chi n^2$$

for some universal  $C_\chi$ . Taylor's theorem with remainder gives

$$f(a_t) \leq \frac{1}{2} \sup_{0 \leq a \leq a_t} |f''(a)| a_t^2 \leq C_\chi n^2 a_t^2,$$

which proves Equation (123). □

**Theorem G.14** (Singular two-point lower bound). *There exists a constant  $c_{\text{sing}} > 0$ , depending only on  $\sigma$ , such that for all sufficiently large  $m$ ,*

$$\mathfrak{R}_{m,n}(\{\eta_0, \eta_{t_{m,n}}\}; d_q) \geq c_{\text{sing}} \frac{\sigma}{m^{1/4} \sqrt{n}}, \quad t_{m,n} := c_0 \frac{\sigma}{m^{1/4} \sqrt{n}}, \quad (124)$$

provided  $c_0 > 0$  is chosen sufficiently small. Consequently,

$$\inf_{\hat{\eta}} \sup_{0 \leq t \leq t_{m,n}} \mathbb{E}_{\eta_t} [d_q(\hat{\eta}, \eta_t)] \gtrsim \frac{1}{m^{1/4} \sqrt{n}}.$$

*Proof.* By Theorem G.2, the reduced experiment is Blackwell equivalent to the full grouped-data experiment, so it suffices to prove the lower bound for the reduced one-task laws  $\tilde{Q}_t^{(n)}$ .

For  $t = t_{m,n}$ , we have

$$a_t = \frac{t^2}{\sigma^2} = \frac{c_0^2}{\sqrt{m}n}.$$

Hence  $a_t \leq 1/(4n)$  for all sufficiently large  $m$ . By Equation (123),

$$\chi^2(\tilde{Q}_t^{(n)}, \tilde{Q}_0^{(n)}) \leq C_\chi \frac{c_0^4}{m}.$$

For the  $m$ -task product experiment,  $\chi^2$  tensorizes:

$$1 + \chi^2((\tilde{Q}_t^{(n)})^{\otimes m}, (\tilde{Q}_0^{(n)})^{\otimes m}) = \left(1 + \chi^2(\tilde{Q}_t^{(n)}, \tilde{Q}_0^{(n)})\right)^m \leq \exp(C_\chi c_0^4).$$

Choose  $c_0 > 0$  small enough that

$$\exp(C_\chi c_0^4) - 1 \leq \frac{1}{2}.$$

Then

$$\text{TV}((\tilde{Q}_t^{(n)})^{\otimes m}, (\tilde{Q}_0^{(n)})^{\otimes m}) \leq \sqrt{\frac{1}{2} \chi^2((\tilde{Q}_t^{(n)})^{\otimes m}, (\tilde{Q}_0^{(n)})^{\otimes m})} \leq \frac{1}{2}.$$

Applying Le Cam's inequality Equation (109) and using Equation (102),

$$\inf_{\hat{\eta}} \sup_{\eta \in \{\eta_0, \eta_t\}} \mathbb{E}_\eta[d_q(\hat{\eta}, \eta)] \geq \frac{d_q(\eta_t, \eta_0)}{4} \left(1 - \frac{1}{2}\right) = \frac{2t}{8} = \frac{t}{4}.$$

Substituting  $t = t_{m,n}$  proves Equation (124).  $\square$

*Remark G.15* (Sharpness relative to the upper bound). The singular lower bound of Theorem G.14 matches the  $\kappa = 2$  structural upper bound from Appendix F up to logarithmic factors, since the posterior rate there is  $O(\delta_{m,n}^{1/2}/\sqrt{n})$  and  $\delta_{m,n} \asymp m^{-1/2}$  up to logs when the model dimensions are fixed.

*Remark G.16* (What Appendix G shows). The lower bounds separate three distinct sources of difficulty:

1. *grammar complexity*: even with a perfectly known module library, the latent task-family law cannot be recovered faster than  $\sqrt{K_{\text{eff}}/m}$ ;
2. *module estimation*: when the benchmark exposes modules in isolation, one inherits the classical  $1/\sqrt{mn}$ -type difficulty of estimating shared operators;
3. *singularity*: in overfitted duplicated-module families, structural estimation slows to  $m^{-1/4}n^{-1/2}$ , even though predictive learning remains much easier.

This is precisely the statistical phase structure that the main paper turns into benchmark hardness labels.

## H. Experimental Protocol

This appendix specifies the benchmark generators, posterior computations, evaluation metrics, and plotting conventions used throughout Section 5. Every experiment is run in the exact reduced experiment  $(X_i, \widehat{W}_i)_{i=1}^m$ , which is Blackwell equivalent to the full grouped-data experiment by Appendix B. All reported quantities therefore correspond directly to the proved model class, rather than to an external approximation layer. Unless otherwise stated, all numerical calculations are performed in double precision, and all figures report medians across 32 independent replications.

### H.1. Benchmark philosophy and released hardness labels

Each benchmark cell is generated from a theorem-matched stochastic task grammar and released with an explicit hardness label tuple

$$h = (K_{\text{eff}}, \Delta_{\text{ctx}}, \lambda_{\text{lin}}, \gamma_{\text{comp}}, \kappa, s_{\text{anc}}).$$

These labels are computed from the definitions in the main text and not fit post hoc.

**Effective grammar complexity.** For every benchmark family,

$$K_{\text{eff}} = |\mathcal{L}| - 1.$$

**Contextual observability.** When relevant, we compute

$$\Delta_{\text{ctx}}(\eta^*) = \min_{a \neq b} \sum_{(u,v) \in \mathcal{C}_{\mathcal{L}}} \nu_{\pi^*}(u,v) \|A_v^*(A_a^* - A_b^*)A_u^*\|_F^2$$

exactly from the benchmark definition.

**Linearized observability.** For the regular families,  $\lambda_{\text{lin}}(\eta^*)$  is computed as the smallest singular value squared of the weighted Jacobian of the program-operator map,

$$\lambda_{\text{lin}}(\eta^*) = \inf_{\|H\|_{\text{fib}}=1} \|D\Phi_{\Theta^*}(H)\|_{\pi^*}^2,$$

using a finite matrix representation in the theorem ambient space.

**Component separation.** We compute

$$\gamma_{\text{comp}}(\eta^*) = \min_{z \neq z'} \|A_z^* - A_{z'}^*\|_F$$

by explicit enumeration of the admissible language.

**Singularity order.** We assign  $\kappa = 1$  for the regular identifiable families and  $\kappa = 2$  for the duplicated-module singular family proved in Appendix F. For the exact non-identifiable observability cell  $\rho = 0$ , we do not assign a finite  $\kappa$ ; instead the cell is labeled as *non-identifiable*.

**Anchor count.** For the anchored families,  $s_{\text{anc}} = s$ . For all other benchmark families,  $s_{\text{anc}} = 0$ .

## H.2. Benchmark suites

We use four theorem-native suites, plus a matched no-anchor comparison family for the anchor experiment.

### H.2.1. REGULAR IDENTIFIABLE SUITE

For  $q \in \{1, 2, 4\}$ , define

$$\mathcal{L}_q^{\text{reg}} = \{(1), (1, 2), \dots, (1, 2, \dots, q+1)\}, \quad r = q+1, \quad d = 2, \quad n = 32.$$

The true grammar law is uniform on  $\mathcal{L}_q^{\text{reg}}$ , and the module library is diagonal:

$$A_a^* = \begin{pmatrix} 0.72 + 0.03a & 0 \\ 0 & 0.35 + 0.02a \end{pmatrix}, \quad a = 1, \dots, q+1.$$

This suite is used for the regular predictive and structural sanity checks shown in Appendix Figures 4 and 5. In the appendix structural figure, the left panel uses  $q \in \{1, 2, 4\}$ , the normalized structural comparison focuses on  $q \in \{1, 2\}$ , and the direct local inverse inset uses the  $q = 4$  cell.

### H.2.2. OBSERVABILITY SUITE

The observability phase-transition family is

$$\mathcal{L}_{\text{obs}} = \{(1, 3), (2, 3)\}, \quad \pi^* = \left(\frac{1}{2}, \frac{1}{2}\right), \quad d = 2, \quad m = 256.$$

The module library is

$$A_1^* = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \quad A_2^* = \begin{pmatrix} 0 & 0 \\ 2 & 0 \end{pmatrix}, \quad A_3^*(\rho) = \begin{pmatrix} 1 & 0 \\ 0 & \rho \end{pmatrix},$$

with

$$\rho \in \{0, 0.05, 0.10, 0.20, 0.40, 0.80\}, \quad n \in \{4, 8, 16, 32, 64\}.$$

For this family,

$$\Delta_{\text{ctx}}(\eta_\rho^*) = \rho^2/2, \quad \gamma_{\text{comp}}(\eta_\rho^*) = \rho.$$

The cell  $\rho = 0$  is exactly non-identifiable and is treated separately in the right panel of Figure 1.

### H.2.3. SINGULAR DUPLICATED-MODULE SUITE

The singular family is

$$\mathcal{L}_{\text{sing}} = \{(1), (2)\}, \quad \pi^* = \left(\frac{1}{2}, \frac{1}{2}\right), \quad d = 1,$$

with

$$A_1(t) = 1 + t, \quad A_2(t) = 1 - t, \quad t^* = 0.$$

We evaluate

$$n \in \{4, 8, 16, 32\}, \quad m \in \{64, 128, 256, 512, 1024\}.$$

In this family,

$$d_q(\eta_t, \eta_0) = 2t, \quad \kappa = 2.$$

The right panel of Figure 2 additionally evaluates the local quartic geometry at

$$t \in \{0.003, 0.0045, 0.006, 0.008, 0.011, 0.015, 0.021, 0.030, 0.042, 0.060\}$$

for each  $n \in \{4, 8, 16, 32\}$ .

### H.2.4. ANCHORED AND MATCHED NO-ANCHOR SUITES

For  $s \in \{1, 2, 4, 8\}$ , define the anchored family

$$\mathcal{L}_s^{\text{anc}} = \{(1), \dots, (s)\} \cup \{(s+1), (s+1, s+2)\},$$

with  $r = s + 2$ ,  $d = 2$ , fixed  $n = 16$ , and grammar weights

$$\pi_{(j)}^* = \frac{\beta}{s}, \quad j = 1, \dots, s, \quad \beta = \frac{1}{2},$$

while the remaining mass is split uniformly over  $(s+1)$  and  $(s+1, s+2)$ . The module library is diagonal:

$$A_a^* = \begin{pmatrix} 0.80 + 0.04a & 0 \\ 0 & 0.25 + 0.03a \end{pmatrix}, \quad a = 1, \dots, s+2.$$

To isolate the value of explicit anchor exposure, we construct a matched no-anchor family

$$\mathcal{L}_s^{\text{noanc}} = \{(1, s+1), \dots, (s, s+1)\} \cup \{(s+1), (s+1, s+2)\},$$

with the same  $r$ , the same module library, and the same grammar weights. Thus the only qualitative difference between the two suites is that in  $\mathcal{L}_s^{\text{anc}}$  the anchor modules appear in singleton tasks, while in  $\mathcal{L}_s^{\text{noanc}}$  they only appear through a shared context.

The right panel of Figure 3 uses the matched comparison at

$$s = 4, \quad n = 16, \quad m \in \{64, 128, 256, 512\}.$$

### H.3. Data generation in the reduced experiment

For every suite and every replication:

1. sample  $m$  latent programs  $Z_1, \dots, Z_m$  independently from the true grammar law  $\pi^*$ ;
2. for each task  $i$ , draw  $x_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d), j = 1, \dots, n$ ;

3. form responses

$$y_{ij} = A_{Z_i}^* x_{ij} + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2 I_d),$$

with  $\sigma = 0.15$ ;

4. compute the exact reduced statistic

$$\widehat{W}_i = Y_i X_i^\top (X_i X_i^\top)^{-1}.$$

All posterior computations are then performed on  $(X_i, \widehat{W}_i)_{i=1}^m$ .

### H.4. Posterior computation

#### H.4.1. REGULAR AND ANCHORED SUITES

For the regular and anchor families, we restrict both the generator and fitted model to the theorem-matched diagonal subfamily. Thus each module is parameterized as

$$A_a = \text{diag}(u_a, v_a).$$

The grammar weights are fixed at the true values in the anchored families and free in the regular family. In all cases, inference uses the *exact reduced likelihood*. Posterior computation proceeds as follows:

1. compute a MAP estimator using deterministic L-BFGS optimization from multiple random initializations;
2. evaluate the exact Hessian of the negative log posterior at the MAP by automatic differentiation;
3. approximate the posterior by the corresponding Gaussian Laplace approximation;
4. compute posterior quantiles from draws of this Gaussian approximation.

The regular suites use 8 optimization restarts and 192 posterior draws. For the  $q = 4$  structural validation cell, we increase the computational budget to 12 restarts and 320 posterior draws. In the anchored suite we use 6 restarts and 160 posterior draws by default, increase to 10 restarts and 240 draws for  $s = 8$ , and to 14 restarts and 384 draws for  $s = 1$ , which is the visibly noisiest case.

#### H.4.2. OBSERVABILITY SUITE

For the observability family, posterior inference is exact on a rectangular grid over

$$(\pi, \rho) \in [\pi_{\min}, \pi_{\max}] \times [\rho_{\min}, \rho_{\max}] = [10^{-3}, 1 - 10^{-3}] \times [0, 0.9].$$

We use 161 grid points in  $\pi$  and 241 grid points in  $\rho$ , together with trapezoidal quadrature weights. The grammar law is parameterized as

$$\pi = (\pi, 1 - \pi),$$

where  $\pi$  denotes the probability of the first program (1, 3). Because the reduced family is exactly two-dimensional, no approximation layer is required.

### H.4.3. SINGULAR SUITE

For the duplicated-module family, posterior inference is exact on a one-dimensional grid

$$t \in [0, 0.2]$$

with 1601 grid points and trapezoidal quadrature weights. The predictive Hellinger geometry is computed exactly by Gauss–Hermite quadrature of order 80. No Laplace approximation or MCMC is used in this suite.

## H.5. Evaluation metrics

**Predictive posterior radius.** For the regular predictive sanity check and the singular-family predictive panel, we use

$$R_{\text{pred},0.9} = \inf\{r : \Pi_m(\bar{h}_X(\eta, \eta^*) \leq r) \geq 0.9\}.$$

In the regular suite,  $\bar{h}_X$  is estimated on a held-out design bank using exact reduced densities and a fixed proposal bank sampled from the true predictive law. In the singular suite,  $\bar{h}_X$  is computed by exact Gauss–Hermite quadrature.

**Structural posterior radius.** In the regular and singular suites, the structural posterior radius is

$$R_{\text{str},0.9} = \inf\{r : \Pi_m(d_q(\eta, \eta^*) \leq r) \geq 0.9\}.$$

For the regular suite,  $d_q$  is evaluated by exact enumeration of the automorphism group of the language. For the duplicated family, this reduces to

$$d_q(\eta_t, \eta_0) = 2t.$$

In the observability suite, the fitted family is two-dimensional, and the restricted quotient metric becomes

$$d_{\text{obs}}((\pi, \rho), (\pi^*, \rho^*)) = 2|\pi - \pi^*| + |\rho - \rho^*|,$$

which equals the general quotient distance restricted to the family with fixed  $A_1$  and  $A_2$ . The right panel of Figure 1 therefore reports

$$R_{\text{str},0.9} = \inf\{r : \Pi_m(d_{\text{obs}} \leq r) \geq 0.9\}.$$

**Latent-program recovery.** For synthetic data the true latent programs  $Z_i^*$  are known, so we report

$$M_Z = \frac{1}{m} \sum_{i=1}^m \Pi_m(Z_i = Z_i^*).$$

In the observability suite this is computed exactly by integrating responsibilities over the exact posterior grid.

**Anchor-module radius.** In the anchor experiments, the structural target is the anchor-only library error

$$R_{\text{anc},0.9} = \inf\left\{r : \Pi_m\left(\min_{\sigma \in S_s} \sum_{j=1}^s \|A_j - A_{\sigma(j)}^*\|_F \leq r\right) \geq 0.9\right\}.$$

Thus the anchor-module error is the *sum* over  $s$  anchor modules after optimal assignment. This is why the left panel of Figure 3 uses the per-anchor normalization

$$\frac{R_{\text{anc},0.9} \sqrt{mn}}{d s^{3/2}},$$

rather than  $R_{\text{anc},0.9}/(d\sqrt{s/(mn)})$ .

## H.6. Figure-specific details

**Appendix Figure 4.** We evaluate the regular suite at

$$q \in \{1, 2, 4\}, \quad m \in \{64, 128, 256, 512, 1024\}, \quad n = 32.$$

The figure reports  $R_{\text{pred},0.9}$  and its normalization by

$$\delta_{m,n} = \left(\frac{K_{\text{eff}} \log m + r d^2 \log(mn)}{m}\right)^{1/2}.$$

**Appendix Figure 5.** The structural regular figure uses the same regular suite and grid. The left panel reports  $R_{\text{str},0.9}$ ; the normalized panel focuses on the classes with the cleanest main-paper comparison ( $q = 1, 2$ ); and the inset directly evaluates local perturbations around the  $q = 4$  truth. Specifically, we draw random zero-sum grammar perturbations and random module perturbation directions, scale them by

$$\varepsilon \in \{0.003, 0.006, 0.009, 0.012, 0.018\},$$

and plot  $\bar{h}_X(\eta, \eta^*)^2$  against  $d_q(\eta, \eta^*)^2$ .

**Figure 1.** For the observability phase transition we evaluate

$$\rho \in \{0, 0.05, 0.10, 0.20, 0.40, 0.80\}, \quad n \in \{4, 8, 16, 32, 64\}, \quad m = 256.$$

The left panel plots  $1 - M_Z$  against

$$n\Delta_{\text{ctx}} = n\rho^2/2,$$

using all positive- $\rho$  cells. The right panel isolates the  $\rho = 0$  family and overlays the corresponding structural and symmetry floors.

**Figure 2.** For the singular family we evaluate

$$n \in \{4, 8, 16, 32\}, \quad m \in \{64, 128, 256, 512, 1024\}.$$

The normalized panels use

$$\frac{R_{\text{pred},0.9}}{\delta_{m,n}} \quad \text{and} \quad \frac{R_{\text{str},0.9}\sqrt{n}}{\delta_{m,n}^{1/2}},$$

with

$$\delta_{m,n} = \left( \frac{\log m + 2 \log(mn)}{m} \right)^{1/2},$$

corresponding to  $K_{\text{eff}} = 1$ ,  $r = 2$ , and  $d = 1$ . The quartic-geometry panel uses the explicit  $(n, t)$ -grid listed above and a slope-1 guide on log-log axes.

**Figure 3.** For the anchored family we evaluate

$$s \in \{1, 2, 4, 8\}, \quad m \in \{64, 128, 256, 512\}, \quad n = 16.$$

The left panel uses the per-anchor normalized quantity

$$R_{\text{anc},0.9}\sqrt{mn}/(d s^{3/2}).$$

The right panel fixes  $s = 4$ ,  $n = 16$ , and compares the anchored family against the matched no-anchor family. Because the  $s = 1$  regime is empirically the noisiest, the anchor figure reports interquartile bands rather than 10th–90th percentile bands.

## H.7. Reproducibility and numerical conventions

Each benchmark cell is rerun with 32 independent random seeds. The plotting scripts use deterministic pseudo-random initialization schedules tied to  $(s, m, \rho, n, \text{family})$  so that the experiments are exactly reproducible. GPU acceleration is used when available, but all scripts fall back to CPU execution without changing the mathematical computations. No figure in the paper relies on a training heuristic, model-selection sweep, or unreported hyperparameter search: every numerical choice used in the final plots is fixed in advance by the suite-specific protocol described above.

## H.8. Criteria for successful theorem validation

The experiments are considered successful if they exhibit the following signatures:

1. **Regular predictive sanity check:** normalization by  $\delta_{m,n}$  substantially stabilizes the predictive posterior radius.

2. **Regular structural sanity check:** normalization by  $\delta_{m,n}$  stabilizes the structural posterior radius in the regular regime, and local perturbations satisfy an approximately linear relation between  $\bar{h}_X^2$  and  $d_q^2$ .
3. **Observability:**  $1 - M_Z$  approximately collapses as a function of  $n\Delta_{\text{ctx}}$ , while the  $\rho = 0$  cell exhibits a visible nonzero uncertainty floor.
4. **Singularity:** predictive radii are regular after  $\delta_{m,n}$ -normalization, but structural radii require the slower singular normalization  $\delta_{m,n}^{1/2}/\sqrt{n}$ .
5. **Anchors:** the per-anchor normalized anchor radius remains order one across the anchor-count grid, and explicit anchor cells consistently outperform matched no-anchor cells.

These signatures are exactly the empirical counterparts of the theorems in Section 3. The benchmark is therefore not only derived from the theory in construction, but validated by checking the theorem-predicted phase structure directly.

## I. Additional Plots

This appendix collects supplementary theorem-validation plots that are useful for completeness but not central to the main-text narrative. In the main paper we focus on the three benchmark regimes that most directly expose our new hardness parameters—observability, singularity, and anchor exposure. The plots below instead serve as *regular-regime sanity checks*: they verify that, when the inverse map is well behaved, the posterior radii behave in the regular way predicted by the theory.

More precisely, Appendix Figure 4 validates the predictive contraction theorem in the regular identifiable family, while Appendix Figure 5 validates the corresponding regular structural behavior and directly visualizes the local inverse law

$$\bar{h}_X(\eta, \eta^*)^2 \asymp d_q(\eta, \eta^*)^2$$

through a local perturbation study. These figures are relegated to the appendix only because they play a supporting role: they show that the regular regime behaves as expected, whereas the main contribution of the paper lies in isolating and benchmarking the nontrivial regimes in which observability vanishes, singularity slows structural learning, or benchmark-cell design changes the effective statistical phase.

**Regular predictive sanity check.** The first supplementary figure reports the regular identifiable family

$$\mathcal{L}_q^{\text{reg}} = \{(1), (1, 2), \dots, (1, 2, \dots, q + 1)\}, \quad q \in \{1, 2, 4\},$$

with  $n = 32$ . The left panel shows the raw predictive posterior radius  $R_{\text{pred},0.9}$ , while the right panel normalizes by the theoretical predictive scale

$$\delta_{m,n} = \left( \frac{K_{\text{eff}} \log m + rd^2 \log(mn)}{m} \right)^{1/2}.$$

The substantive point is not exact overlap across grammar classes, which would require class-independent constants, but rather that the normalization substantially stabilizes the family of curves and removes most of the  $m$ -dependence.

**Regular structural sanity check and local inverse validation.** The second supplementary figure uses the same regular family. Its left panel reports the structural posterior radius  $R_{\text{str},0.9}$ . The normalized panel focuses on the lower-complexity classes  $q \in \{1, 2\}$ , where the rate-normalized comparison is the clearest visually. The more important diagnostic is the inset, which directly probes the regular local inverse theorem by generating small perturbations around the  $q = 4$  truth and plotting

$$\bar{h}_X(\eta, \eta^*)^2 \quad \text{against} \quad d_q(\eta, \eta^*)^2.$$

The approximate linearity of this relation is exactly what one expects when the inverse exponent is  $\kappa = 1$ . Thus the inset is the local geometric counterpart to the regular structural-contraction theorem: the posterior rate behaves regularly because the predictive geometry is locally Euclidean in the quotient parameter.

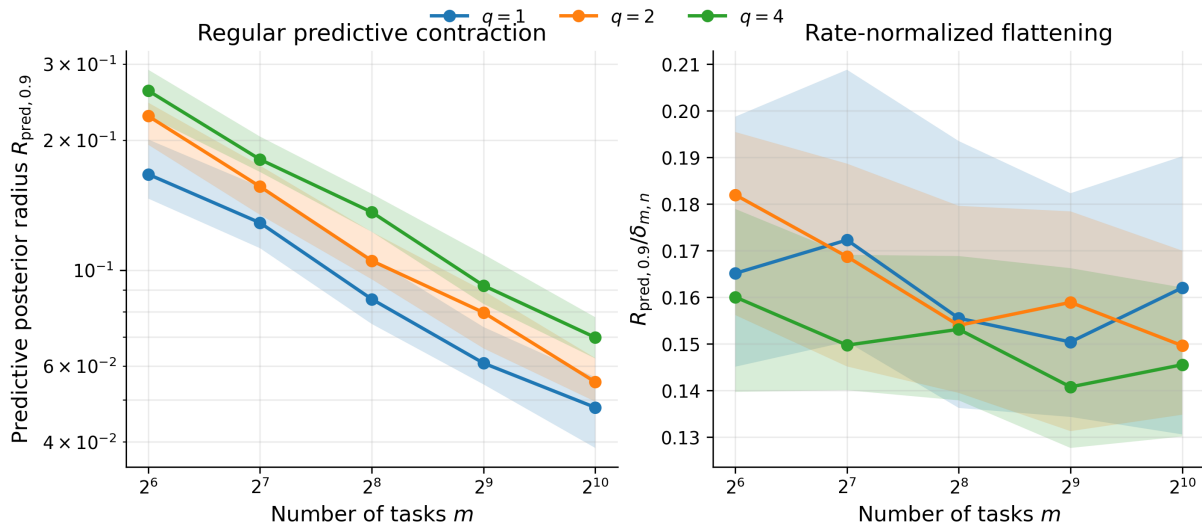
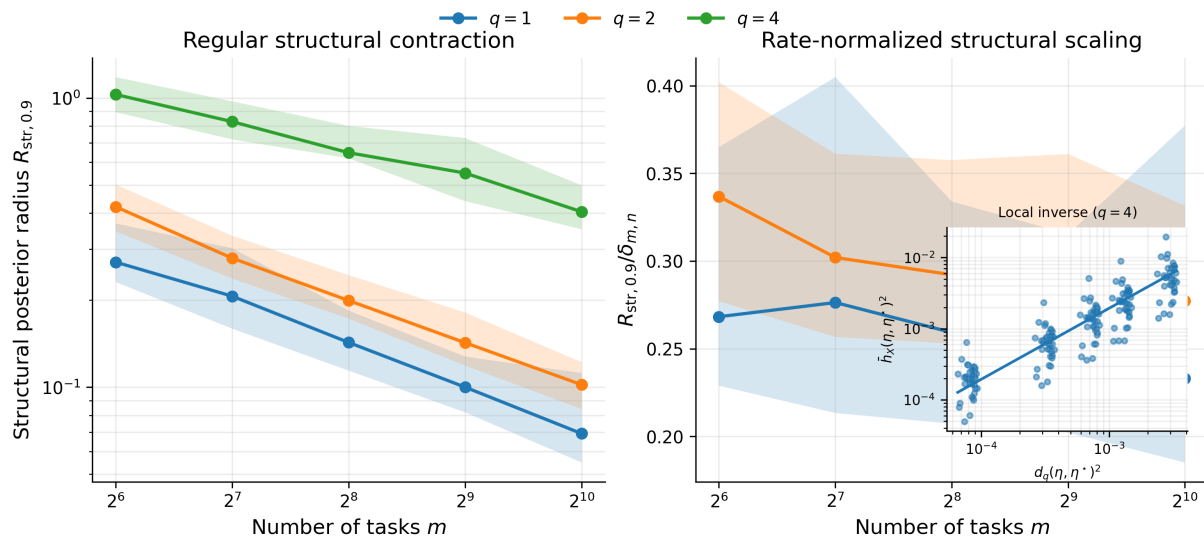


Figure 4. **Regular predictive contraction in the stochastic linear task-grammar family.** *Left:* the posterior predictive radius  $R_{\text{pred},0.9}$  decreases with the number of tasks  $m$ , with larger grammar complexity  $q$  yielding larger radii. *Right:* dividing by the theoretical predictive scale  $\delta_{m,n}$  substantially flattens the curves, consistent with the predicted regular predictive rate up to model-class-dependent constants. Curves show medians across replications; shaded bands indicate 10th–90th percentile ranges.

**How these plots complement the main text.** Taken together, Appendix Figures 4 and 5 show that the regular regime behaves exactly as the theory says it should: predictive and structural uncertainty contract at the same order, and the local inverse map is first-order nondegenerate. The main text then uses this regular baseline as the reference point against which the more novel benchmark regimes are contrasted. In particular, Figure 1 shows what happens when observability collapses, Figure 2 shows what happens when the inverse map becomes singular, and Figure 3 shows how benchmark design can move the problem between qualitatively different statistical regimes.



**Figure 5. Regular structural contraction and direct local inverse validation.** *Left:* the structural posterior radius  $R_{\text{str},0.9}$  decreases with the number of tasks  $m$ , with larger grammar classes yielding larger radii. *Right:* normalization by the regular structural scale  $\delta_{m,n}$  stabilizes the lower-complexity regular classes. The inset directly validates the regular local inverse law, showing an approximately linear relationship between  $\bar{h}_X(\eta, \eta^*)^2$  and  $d_q(\eta, \eta^*)^2$  in the  $q = 4$  regular regime. Curves show medians across replications; shaded bands indicate 10th–90th percentile ranges.