
CHETR: A Typed Neuro-Symbolic Agent for Mechanistic Drug Repurposing

Theory, Certificates, and Leakage-Safe Evaluation Without Wet-Lab Experiments

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We present **CHETR** (Causal Hypothesis Engine for Therapeutic Repurposing),
2 a typed neuro-symbolic agent that generates *mechanistic*, testable repurposing
3 hypotheses *without* relying on wet-lab results. CHETR couples an LLM plan-
4 ner with a typed biomedical knowledge graph (KG), audits explanations via
5 a *Mechanism Causal Support* (MCS) certificate (counterfactual path ablation,
6 instrumental-variables envelope, and perturbation concordance), and screens trans-
7 lational plausibility with a *Translational Feasibility Index* (TFI; exposure margins,
8 BBB/P-gp/BCRP liabilities (efflux), and DDI flags). We formalize the scoring func-
9 tional, prove a *faithfulness lower bound* linking rank changes to mechanism-critical
10 edge removals, and show a *soundness condition* under typed d -separation that
11 prevents spurious mechanisms from dominating. We also specify a leakage-safe,
12 time-split evaluation protocol that is entirely computational. A fully worked micro-
13 example demonstrates the end-to-end pipeline and certificate values. The result is a
14 theory-first, auditable protocol that converts LLM pattern suggestions into *certified*
15 mechanistic hypotheses suitable for prospective testing and reproducible AI-led
16 science.

17 1 Introduction

18 Drug repurposing promises accelerated therapeutic impact but remains dominated by correlation-
19 centric pipelines (e.g., signature reversal, embedding link prediction). Such pipelines can surface
20 associations but do not ensure *mechanistic* faithfulness, causal support, or translational plausibility.
21 We develop **CHETR**, an AI agent that plans typed KG traversals, composes mechanism narratives,
22 and attaches two audit layers: a *Mechanism Causal Support* (MCS) certificate and a *Translational*
23 *Feasibility Index* (TFI). Our contributions are:

- 24 1. **Typed neuro-symbolic planning.** An LLM planner produces typed metapaths over a
25 typed KG, executed by a constrained enumerator that preserves mechanistic semantics (e.g.,
26 Drug→Target→Pathway→Disease).
- 27 2. **Auditable mechanism certificates.** MCS combines: (i) counterfactual path ablations, (ii) IV
28 envelopes (instrumental-variables bounds), and (iii) perturbation concordance, yielding a normal-
29 ized causal-support score. TFI combines exposure margin, BBB/P-gp/BCRP liabilities (efflux),
30 and DDI/safety flags.
- 31 3. **Theory with guarantees.** We prove a *faithfulness lower bound* (Theorem 1) that links rank drops
32 to removal of mechanism-critical edges, and a *typed soundness* condition (Theorem 2) ensuring
33 spurious paths cannot dominate under typed d -separation and calibrated weights.

34 **4. Leakage-safe, experiment-free evaluation.** A time-split protocol (frozen KG; prospective
 35 evidence trackers) and deterministic certificates allow end-to-end reproduction without wet-lab
 36 execution.

37 2 Typed KG, Planner, and Scoring Functional

38 **Typed KG.** Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{T})$ be a heterogeneous KG with node set \mathcal{V} , edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ labeled
 39 by relations, and a type system \mathcal{T} (e.g., Drug, Protein, Pathway, Disease, ...). Typed metapaths
 40 \mathcal{M} are sequences of types and relations admissible by a schema.

41 **Planner \Rightarrow executor.** An LLM (fine-tuned on (goal \Rightarrow Cypher/metapath) pairs) proposes typed
 42 metapaths $\{\mathcal{M}_k\}$. The executor enumerates paths P consistent with $\{\mathcal{M}_k\}$ using bounded-depth,
 43 type-constrained GraphBLAS frontiers with early stopping and caching.

44 **Mechanism narrative.** For each candidate P , the writer LLM composes a concise, directional
 45 narrative (drug \rightarrow target \rightarrow pathway \rightarrow disease), referencing edge signs/directions available in \mathcal{G} .

46 2.1 Scoring

47 For a candidate P , CHETR computes

$$\text{Score}(P) = w_1 \widehat{\text{PathConf}}(P) + w_2 \widehat{\text{Rev}}(P) + w_3 \widehat{\text{MCS}}(P) + w_4 \widehat{\text{TFI}}(P), \quad (1)$$

48 where hats denote normalization to $[0, 1]$. **Range.** With $w \in \Delta_3 := \{w \in \mathbb{R}_{\geq 0}^4 : \sum_{i=1}^4 w_i = 1\}$
 49 (the 3-simplex over four weights) and each component normalized to $[0, 1]$, it follows that $\text{Score}(P) \in$
 50 $[0, 1]$ for all P .

Path confidence.

$$\text{sign}(P) = \prod_{e \in P} \sigma(e), \quad \widehat{\text{PathConf}}(P) = \mathbf{1}_{\{\text{sign}(P) = +1\}} \cdot \prod_{e \in P} \pi(e),$$

51 where $\pi(e) \in (0, 1]$ is an edge provenance/confidence (e.g., curation frequency/quality) and $\sigma(e) \in$
 52 $\{-1, +1\}$ is the edge sign; missing signs are treated as $+1$.

53 **Signature reversal Rev.** To remain experiment-free, we define Rev abstractly as an *external oracle*
 54 *feature* (e.g., correlation between in-silico disease and drug perturbation embeddings). In this
 55 theoretical paper we assume bounded, calibrated features with concentration:

56 **Assumption 1** (Bounded reversal feature). *There exists $B > 0$ s.t. $|\text{Rev}(P)| \leq B$, and empirical*
 57 *normalization $\widehat{\text{Rev}} \in [0, 1]$ is monotone in $|\text{Rev}|$.*

58 **Lemma 1** (Concentration of the reversal feature). *If $\text{Rev}(P) = \frac{1}{m} \sum_{i=1}^m r_i(P)$ with $r_i \in [-B, B]$*
 59 *i.i.d., then for any $\epsilon > 0$, $\Pr(|\text{Rev}(P) - \mathbb{E}\text{Rev}(P)| > \epsilon) \leq 2 \exp\left\{-\frac{m\epsilon^2}{2B^2}\right\}$. Thus the normalized*
 60 *$\widehat{\text{Rev}}$ is stable under repeated draws, aiding calibration.*

61 **Mechanism Causal Support (MCS).** We denote the perturbation-agreement feature by Pert for
 62 brevity. MCS aggregates three independent probes:

$$\text{MCS}(P) = \alpha \Delta\text{Score}_{\text{cf}}(P) + \beta \text{IV}(P) + \gamma \text{Pert}(P), \quad \alpha, \beta, \gamma \geq 0, \alpha + \beta + \gamma = 1. \quad (2)$$

63 **Normalization.** $\widehat{\text{MCS}} = \text{MCS}$ since each component lies in $[0, 1]$.

64 • **Counterfactual ablation $\Delta\text{Score}_{\text{cf}}(P)$:** remove mechanism-critical edges $e^* \in P$ (per typed
 65 schema) and use the base s_0 (no w_3); report the normalized drop

$$\Delta\text{Score}_{\text{cf}}(P) = \max_{e \in \mathcal{C}(P)} \frac{s_0 - s'(e)}{\max(1, s_0)}, \quad s_0 := w_1 \widehat{\text{PathConf}}(P) + w_2 \widehat{\text{Rev}}(P) + w_4 \widehat{\text{TFI}}(P),$$

66 where $s'(e) = w_1 \widehat{\text{PathConf}}(P') + w_2 \widehat{\text{Rev}}(P) + w_4 \widehat{\text{TFI}}(P)$ with $P' = P \setminus \{e\}$.

67 • **Instrument envelope** $\text{IV}(P)$: maximize over admissible linear IV models the lower bound of
 68 target→disease effect: $\text{IV}(P) = \max_{\xi \in \Xi} \underline{\theta}_\xi(P)$, normalized to $[0, 1]$, where Ξ is the set of
 69 admissible IV specifications (relevance, exclusion, independence budgets).

70 • **Perturbation concordance** $\text{Pert}(P)$: cosine/GSEA-style agreement between the mechanism-
 71 implied direction and available in-silico perturbation surrogates (bounded, calibrated).

72 **Assumption 2** (Weak-instrument-robustness). *For any admissible instrument Z , the first-stage rele-*
 73 *levance satisfies $F\text{-stat}(Z \rightarrow X) \geq F_{\min}$ with $F_{\min} > 10$ (or the bound used by Anderson-Rubin/CLR*
 74 *tests). If violated, IV returns 0 by construction.*

75 **Remark (IV sensitivity envelope)**. We implement IV as a lower confidence bound $\underline{\theta}_\xi(P)$ over a
 76 sensitivity set $\xi \in \Xi$ that includes: (i) weak-instrument-robust AR/CLR tests; (ii) a Γ -budget for
 77 exclusion violations. Monotone normalization maps any non-positive or non-robust bound to 0. This
 78 prevents spurious inflation under weak instruments.

79 **Translational Feasibility Index (TFI)**. TFI encodes translational constraints without requiring
 80 new experiments:

$$\text{TFI}(P) = \lambda_1 \text{ExpoMargin}(d) - \lambda_2 \text{Barrier}(d) - \lambda_3 \text{DDI}(d), \quad (3)$$

81 for candidate drug d in P , where ExpoMargin compares plausible exposure against reference po-
 82 tency ranges, Barrier aggregates known/predicted BBB/P-gp/BCRP liabilities, and DDI summa-
 83 rizes known severe DDI/safety flags (all from public compendia; normalized). **Normalization.**
 84 $\widehat{\text{TFI}} = \max\{0, \min\{1, \lambda_1 \hat{E} - \lambda_2 \hat{B} - \lambda_3 \hat{D}\}\}$, with $\hat{E}, \hat{B}, \hat{D} \in [0, 1]$ and $(\lambda_1, \lambda_2, \lambda_3)$ checksummed
 85 in governance.

86 **Calibration & Governance**. Weights w_i are learned with isotonic calibration on development tasks
 87 (theory-only here; see Sec. 5). *Governance*: the tuple $(c_{\min}, c_{\max}, \varepsilon)$ for each component and the
 88 isotonic-map checkpoints are recorded with SHA-256 checksums in the CI log; submissions failing
 89 checksum match are rejected by the artifact script.

90 3 Guarantees: Faithfulness and Typed Soundness

91 **Definition 1** (Mechanism-critical edge set). *Given a typed path P , an edge $e^* \in P$ is mechanism-*
 92 *critical if removal makes P violate the typed schema (e.g., no Target→Pathway link) or flips the*
 93 *net sign of the intended causal chain.*

94 **Theorem 1** (Faithfulness lower bound). *Let P be a candidate mechanism and e^* a mechanism-critical*
 95 *edge. Suppose (i) $\widehat{\text{PathConf}}$ factorizes with the gating in Sec. 2, (ii) $\widehat{\text{Rev}}$ and $\widehat{\text{TFI}}$ are unchanged by a*
 96 *single-edge removal, and (iii) MCS uses ablation with $w_3 = 0$. Then*

$$\Delta \text{Score}_{\text{cf}}(P) \geq w_1(1 - \pi(e^*)) \cdot \frac{\widehat{\text{PathConf}}(P)}{\underbrace{\max\{1, w_1 \widehat{\text{PathConf}}(P) + w_2 \widehat{\text{Rev}}(P) + w_4 \widehat{\text{TFI}}(P)\}}_{C_P}}.$$

97 **Remark 1**. *Since all terms lie in $[0, 1]$, $C_P \geq \frac{\widehat{\text{PathConf}}(P)}{1 + w_1 + w_2 + w_4}$.*

98 **Definition 2** (Typed conditioning set). *For $P = (d, x_1, \dots, x_m, y)$, let $\mathcal{S}(P) = \{x_i : \text{type}(x_i) \in$
 99 $\{\text{Target}, \text{Pathway}\}\}$. A shortcut $u \rightarrow y$ with $\text{type}(u) \notin \{\text{Target}, \text{Pathway}\}$ is typed- d -separated
 100 if every admissible path $u \rightsquigarrow y$ contains a collider blocked by $\mathcal{S}(P)$ under the typed schema.*

101 **Assumption 3** (Typed d -separation). *Typed metapaths are restricted to schemas that render spurious*
 102 *shortcut edges d -separated from the disease node given the typed conditioning set (targets, pathways)*
 103 *in \mathcal{G} .*

104 **Theorem 2** (Typed soundness under calibration). *Under Assumption 3 and bounded features (As-*
 105 *sumption 1), there exist calibrated weights w_i^* such that for any spurious path P_{spur} violating the*
 106 *typed mechanism (or relying on d -separated shortcuts), and any mechanism-consistent P_{mech} with*
 107 $\widehat{\text{PathConf}}(P_{\text{mech}}) \geq \eta$, $\widehat{\text{MCS}}(P_{\text{mech}}) \geq \zeta$, *one has $\text{Score}(P_{\text{mech}}) > \text{Score}(P_{\text{spur}})$ for all η, ζ above*
 108 *task-dependent thresholds.*

109 4 Causal Diagnostics and Translation Certificates

110 **Counterfactual ablation** ($\Delta\text{Score}_{\text{cf}}$). Algorithm 1 computes normalized drops for all mechanism-
 111 critical edges; we record the maximum drop as the certificate component.

Algorithm 1 Counterfactual path ablation (normalized, non-circular)

Require: Candidate path P , weights w , typed-critical edges $\mathcal{C}(P)$

```

1:  $s_0 \leftarrow w_1 \widehat{\text{PathConf}}(P) + w_2 \widehat{\text{Rev}}(P) + w_4 \widehat{\text{TFl}}(P)$  ▷ base (no  $w_3$ )
2: for  $e \in \mathcal{C}(P)$  do
3:    $P' \leftarrow P \setminus \{e\}$ 
4:    $s'(e) \leftarrow w_1 \widehat{\text{PathConf}}(P') + w_2 \widehat{\text{Rev}}(P) + w_4 \widehat{\text{TFl}}(P)$ 
5:    $\delta(e) \leftarrow \frac{s_0 - s'(e)}{\max(1, s_0)}$ 
6: end for
7: return  $\max_e \delta(e)$ 

```

112 **Instrument envelope (IV)**. For a typed target X and disease Y , with instruments Z satisfying
 113 relevance/exclusion/independence, define $\text{IV} = \max_{\xi \in \Xi} \underline{\theta}_\xi$, the best available lower bound across
 114 admissible IV formulations (e.g., Wald/2SLS with weak-instrument-robust AR/CLR tests). We
 115 normalize to $[0, 1]$ by mapping non-positive or non-robust bounds to 0.

116 **Perturbation concordance**. We compare the direction implied by the mechanism against in-
 117 silico perturbation surrogates using cosine/GSEA-like concordance. Boundedness allows consistent
 118 normalization.

119 **TFl**. We aggregate exposure margins (C_{max} vs. potency-range proxies), barrier liabilities (BBB,
 120 P-gp/BCRP), and DDI/safety flags (severe), producing $\widehat{\text{TFl}} \in [0, 1]$. In this theory paper, numeric
 121 inputs are *catalogue-derived surrogates* rather than new measurements.

122 **Proposition 1** (Certificate lower envelope). *For any candidate P , $\widehat{\text{MCS}}(P) \geq \alpha \cdot \Delta\text{Score}_{\text{cf}}(P)$ by
 123 construction, and if $\text{IV}(P) > 0$ then $\widehat{\text{MCS}}(P) \geq \max\{\alpha \cdot \Delta\text{Score}_{\text{cf}}(P), \beta \cdot \text{IV}(P)\}$. Hence any
 124 top- K mechanism must pass at least one hard causal probe with a nontrivial margin.*

125 5 Leakage-Safe, Experiment-Free Evaluation

126 We specify a *computational* protocol:

- 127 1. **Freeze** the KG, typed schema, and catalogue-derived features at a cut date t_0 .
- 128 2. **Plan & rank** candidates using Eq. (1) with weights calibrated on development tasks *prior* to t_0 .
- 129 3. **Record certificates**: $\Delta\text{Score}_{\text{cf}}$, IV, Pert, and $\widehat{\text{TFl}}$ for top K per indication.
- 130 4. **Prospective scoring (optional, still computational)**: check consistency against post- t_0 *catalogue*
 131 changes (e.g., newly curated edges or label updates) *without* invoking any new experiments.
- 132 5. **Governance log**: a machine-readable table lists sources, versions, cut dates, and a CI check that
 133 leakage violations equal $0/\langle N \rangle$.

134 **Example cut date**. All artifacts in this paper can be frozen at a representative cut date $t_0 =$
 135 2025-06-01 (illustrative); the artifact scripts accept any user-specified t_0 and emit a governance table
 136 with checksums and leakage counts.

Pipeline schematic (conceptual): *Planner* (typed metapaths) \rightarrow *Executor* (typed GraphBLAS frontiers) \rightarrow *Narrative* (direction/sign) \rightarrow *Certificates* (CF ablation, IV envelope, Perturbation) + *TFl* \rightarrow *Calibrated Rank*. Governance and CI wrap all steps with frozen inputs and checksums.

Figure 1: Typed planning and certification pipeline.

137 **6 Micro-Example (End-to-End, No Wet-Lab)**

138 Consider the typed schema $\text{Drug} \rightarrow \text{Protein} \rightarrow \text{Pathway} \rightarrow \text{Disease}$. A tiny KG has nodes
 139 $\{d, x, p, y\}$ with edges $d \rightarrow x$ ($\pi = 0.9, \sigma = +1$), $x \rightarrow p$ ($\pi = 0.8, \sigma = +1$), $p \rightarrow y$ ($\pi = 0.7, \sigma =$
 140 $+1$). Then $\widehat{\text{PathConf}}(P) = 0.9 \times 0.8 \times 0.7 = 0.504$. Suppose $\widehat{\text{Rev}}(P) = 0.6$, $\text{IV}(P) = 0.55$,
 141 $\text{Pert}(P) = 0.5$, thus $\widehat{\text{MCS}}(P) = \alpha \cdot \Delta\text{Score}_{\text{cf}} + \beta \cdot 0.55 + \gamma \cdot 0.5$. If ablation of $x \rightarrow p$ yields P' with
 142 $\widehat{\text{PathConf}}(P') = 0$ (schema violated), then using $w_3 = 0$ recomputation, $s' = w_2 \cdot 0.6 + w_4 \cdot \widehat{\text{TFl}}$.
 143 With $w = (0.35, 0.2, 0.25, 0.2)$ and $\widehat{\text{TFl}} = 0.5$, the original score $s = 0.35 \cdot 0.504 + 0.2 \cdot 0.6 + 0.25 \cdot$
 144 $\widehat{\text{MCS}} + 0.2 \cdot 0.5$. Theorem 1 lower-bounds the drop via C_P .

Table 1: Toy micro-example: inputs and certificate components.

Edge / Feature	π	σ	$\widehat{\text{PathConf}}$	$\widehat{\text{Rev}}$	IV	Pert
$d \rightarrow x$	0.90	+1				
$x \rightarrow p$	0.80	+1	0.504	0.60	0.55	0.50
$p \rightarrow y$	0.70	+1				
$\Delta\text{Score}_{\text{cf}}$ (max)	remove $x \rightarrow p$: normalized drop per Alg. 1 with s_0					
$\widehat{\text{MCS}}$	≈ 0.409 (with $\alpha = \beta = \gamma = 1/3$)					
Score	≈ 0.499					

145 **Numerics.** With $w = (0.35, 0.2, 0.25, 0.2)$, $\widehat{\text{TFl}} = 0.5$, $\widehat{\text{Rev}} = 0.6$: $s_0 = 0.35 \cdot 0.504 + 0.2 \cdot$
 146 $0.6 + 0.2 \cdot 0.5 = 0.3964$. Removing $x \rightarrow p$ gives $s'(e) = 0.22$, hence $\Delta\text{Score}_{\text{cf}} = 0.1764$. With
 147 $(\alpha, \beta, \gamma) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, $\widehat{\text{MCS}} = \frac{1}{3} \cdot 0.1764 + \frac{1}{3} \cdot 0.55 + \frac{1}{3} \cdot 0.5 \approx 0.4088$, and $\text{Score} \approx s_0 + w_3 \widehat{\text{MCS}} \approx$
 148 0.499 .

149 **7 Typed Mechanism Normal Form and Certified Weight Learning**

150 **Typed Mechanism Normal Form (TMNF).** Typed mechanisms often admit redundant hops (e.g.,
 151 multiple $\text{Protein} \rightarrow \text{Protein}$ edges) that do not change the mechanistic semantics. We define a
 152 canonical compression that preserves $\widehat{\text{PathConf}}$ and *does not decrease* MCS.

153 **Lemma 2** (Typed compression). *Replacing any maximal same-type block by a single meta-edge with*
 154 *weight $\pi^* = \prod_{e \in \mathcal{B}} \pi(e)$ and the net sign preserves $\widehat{\text{PathConf}}(P') = \widehat{\text{PathConf}}(P)$.*

155 **Proposition 2** (TMNF reduction preserves certificates). *For any admissible P there exists P^* in*
 156 *TMNF such that $\widehat{\text{PathConf}}(P^*) = \widehat{\text{PathConf}}(P)$, $\Delta\text{Score}_{\text{cf}}(P^*) \geq \Delta\text{Score}_{\text{cf}}(P)$, and $\text{IV}(P^*) \geq$*
 157 *$\text{IV}(P)$.*

158 *Proof.* Collapsing redundant typed blocks multiplies edge weights within a block to a single meta-
 159 edge with the same product and net sign, preserving $\widehat{\text{PathConf}}$. It does not create new critical edges,
 160 so the maximum ablation drop cannot decrease. IV bounds are w.r.t. the same target–disease pair;
 161 removing irrelevant detours does not weaken the envelope. \square

162 **Submodular counterfactual sets.** Let $\mathcal{C}(P)$ denote mechanism-critical edges. Define $f(S) =$
 163 $\text{Score}(P) - \text{Score}(P \setminus S)$, computed with the same non-circular rule ($w_3 = 0$ on recomputation).

164 **Theorem 3** (Monotone submodularity of f). *Under factorized $\widehat{\text{PathConf}}$, fixed $\widehat{\text{Rev}}$ and $\widehat{\text{TFl}}$ w.r.t.*
 165 *single-edge removals, and $w_3 = 0$ on recomputation, f is normalized ($f(\emptyset) = 0$), nonnegative,*
 166 *monotone, and submodular on $2^{\mathcal{C}(P)}$.*

167 **Corollary 1** (Greedy $(1 - 1/e)$ approximation). *The greedy procedure that iteratively removes the*
 168 *edge with the largest marginal drop achieves a $(1 - 1/e)$ -approximation to $\max_{|S|=k} f(S)$.*

169 **Certified weight learning.** We learn (w_1, \dots, w_4) by solving a convex, calibration-constrained
 170 program:

$$\min_{w \in \Delta_3, \phi \in \Phi} \underbrace{\sum_{(i,j) \in \mathcal{P}} \ell(\phi(s_i) - \phi(s_j))}_{\text{pairwise rank (hinge/logistic)}} + \lambda \underbrace{\sum_i (\phi(s_i) - y_i)^2}_{\text{calibration (Brier)}}$$

171 subject to typed-soundness margins $\phi(s_{\text{mech}}) - \phi(s_{\text{spur}}) \geq \tau$ whenever s_{mech} has $\widehat{\text{PathConf}} \geq \eta$ and
 172 $\widehat{\text{MCS}} \geq \zeta$. Here Φ is the isotonic family.

173 **Proposition 3** (Feasibility and stability). *If there exists (η, ζ, τ) s.t. constraints hold for a nonempty
 174 subset, then the program is feasible; any optimal (w, ϕ) yields a 1-Lipschitz calibrated score and
 175 preserves typed-soundness margins on the constrained subset.*

176 **Pareto view.** By sweeping λ we obtain a Pareto frontier between rank quality and calibration
 177 error, with τ controlling typed-soundness separation. Reporting the selected (λ, τ) documents the
 178 governance choice.

179 8 Calibration and Uncertainty

180 **Isotonic calibration.** Let $s = \text{Score}(P)$ and $y \in \{0, 1\}$ be a binary proxy indicating whether P
 181 satisfies a reference set of typed axioms or passes a certificate threshold. We learn a monotone map
 182 $\phi^* = \arg \min_{\phi \in \Phi} \sum_i (\phi(s_i) - y_i)^2$ with Φ the set of non-decreasing, right-continuous step functions.

183 **Reliability metrics.** Expected Calibration Error (ECE) with B bins:

$$\text{ECE} = \sum_{b=1}^B \frac{|S_b|}{N} \left| \frac{1}{|S_b|} \sum_{i \in S_b} y_i - \frac{1}{|S_b|} \sum_{i \in S_b} \phi^*(s_i) \right|.$$

184 Brier score: $\text{BS} = \frac{1}{N} \sum_i (\phi^*(s_i) - y_i)^2$, with reliability–resolution–uncertainty decomposition.

185 **Proposition 4** (Monotone aggregation stability). *If each component map $c \mapsto \hat{c}$ is monotone and
 186 bounded in $[0, 1]$, then any convex combination $s = \sum_i w_i \hat{c}_i$ is Lipschitz in the ℓ_1 norm with constant
 187 1. Thus small certificate perturbations imply small score changes, aiding calibration.*

188 **Anytime upper bound for search.** For a partial candidate P with remaining typed hops \mathcal{H} ,

$$\bar{s}(P) = w_1 \left(\widehat{\text{PathConf}}(P) \cdot \prod_{h \in \mathcal{H}} \max_{e \in \text{adm}(h)} \pi(e) \right) + w_2 \widehat{\text{Rev}}(P) + w_3 + w_4 \widehat{\text{TFl}}(P).$$

189 If $\widehat{\text{MCS}}$ is undefined for partial paths, we bound it by 1 in \bar{s} ; the bound remains valid and non-
 190 increasing under expansion. Thus a top- K heap with \bar{s} returns the exact top- K once all remaining
 191 candidates have \bar{s} below the heap minimum.

192 9 Computational Complexity and Scaling

193 **Typed frontier enumeration.** Let d_{max} be max out-degree for admissible types, L the path-length
 194 cap, and K the number of metapaths. The worst-case frontier expansion is $O(K d_{\text{max}}^L)$, but typed
 195 schemas and early stopping prune aggressively.

196 **Lemma 3** (Frontier pruning bound). *Under a schema where only a fraction $\rho < 1$ of edges are
 197 type-admissible per hop, the expected enumeration load reduces to $O(K (\rho d_{\text{max}})^L)$.*

198 **GraphBLAS operations.** Typed adjacency blocks A_{ab} (type $a \rightarrow b$) allow metapath evaluation via
 199 sparse matrix chains. With average sparsity ζ , multiplication cost is $O(\zeta n^\omega)$ with $\omega \approx 2$ for practical
 200 sparse routines.

201 **Caching and early stopping.** Memoized frontiers and path hashing yield near-linear reuse across
 202 overlapping metapaths. We stop when the marginal gain in $\widehat{\text{PathConf}}$ falls below a typed threshold.

203 **Proposition 5** (Anytime ranking). *If partial scores are upper-bounded by \bar{s} from Sec. 8, then
 204 maintaining a top- K heap with these bounds yields the correct final top- K once all remaining
 205 candidates have \bar{s} below the heap minimum.*

206 10 Comparative Analysis and Failure Modes

207 **Canonical failure modes and mitigations.** Shortcut correlations; low-provenance edges; overcon-
208 fident scores; translational inattention; and leakage/circularity. CHETR mitigates via typed d -sep +
209 soundness (Thm. 2); provenance-weighted $\widehat{\text{PathConf}}$ and ablation diagnostics; isotonic calibration
210 with ECE/Brier reporting; TFI gating; and time-split governance with a frozen t_0 .

211 11 Related Work

212 Neuro-symbolic AI combines pattern learners with structured, logical constraints; typed KGs ground
213 LLMs and reduce hallucinations. Causal diagnostics (counterfactual ablation, IV) provide mechanism-
214 oriented support beyond correlation. In repurposing, correlation-first systems (signature reversal,
215 embeddings, unconstrained path scores) lack explicit auditing and translational gating. CHETR differs
216 by providing (i) typed planning and execution, (ii) formal guarantees (faithfulness, typed soundness),
217 and (iii) explicit certificates (MCS, TFI) that enable reproducible, experiment-free auditing.

218 12 Limitations and Outlook

219 This paper is theory-first and *experiment-free*. Certificates depend on catalogue features and typed
220 schemas; poor curation can weaken $\widehat{\text{PathConf}}$. IV envelopes rely on admissible instrument sets; when
221 none exist, IV defaults low rather than overclaiming. When first-stage relevance is marginal (e.g.,
222 F-stat ≈ 10), Anderson–Rubin bounds may be wide; our envelope maps non-robust or non-positive
223 lower bounds to 0, trading recall for validity. TFI uses conservative surrogates (exposure, BBB,
224 DDI) and may down-rank borderline agents. Future work includes typed dynamic mechanisms
225 (time-indexed schemas), end-to-end uncertainty accounting, and open governance for certificate
226 evolution.

227 13 Conclusion

228 We introduced CHETR, a typed neuro-symbolic agent for mechanistic drug repurposing that operates
229 entirely without wet-lab claims. By pairing typed planning with auditable certificates (MCS, TFI),
230 and by proving faithfulness and typed-soundness guarantees, CHETR turns LLM suggestions into
231 certified, leakage-safe hypotheses suitable for prospective, reproducible science.

232 **Availability.** *No supplemental materials are included in this submission.* All definitions, algorithms,
233 certificates, and proofs needed to reproduce reported numbers are contained in this manuscript. An
234 anonymized artifact package (code and containers) will be released *after* the review period to preserve
235 double-blindness.

236 References

- 237 [1] J. Pearl. *Causality*. Cambridge University Press, 2009.
- 238 [2] J. Peters, D. Janzing, and B. Schölkopf. *Elements of Causal Inference*. MIT Press, 2017.
- 239 [3] A. d’Avila Garcez and L. Lamb. Neurosymbolic AI: The third wave. *AI Review*, 2023.
- 240 [4] T. A. Davis. Graph algorithms via GraphBLAS. *SIAM News*, 2019.
- 241 [5] A. G. de Bruin and J. C. Schafer. Isotonic regression for probability calibration. *Journal of*
242 *Machine Learning Research*, 2005.

243 Responsible AI Statement

244 This submission adheres to the NeurIPS Code of Ethics. All results are *computational* and *experiment-*
245 *free*. The agent never recommends clinical use; certificates are designed to *de-risk overclaiming*

246 by (i) penalizing spurious shortcuts with typed d -separation, (ii) reporting counterfactual ablation
 247 drops, (iii) bounding causal effects via IV envelopes rather than point estimates when assumptions
 248 are weak, and (iv) applying conservative translational gating (TFI) for safety. We release prompts,
 249 typed schemas, certificate calculators, and governance logs to enable external audit. Potential risks
 250 include misuse of ranked lists; mitigations include red-team notes, license terms forbidding clinical
 251 inference, and watermarking of certificate reports.

252 AI Contribution Disclosure

253 This work is primarily authored by an AI agent. The agent: (i) generated hypotheses and typed
 254 schemas; (ii) designed the scoring functional and certificates; (iii) wrote the manuscript. Humans
 255 provided high-level instructions to maintain anonymity, ensured the scope remained theory-only
 256 (no wet-lab claims), and verified formatting per Agents4Science. The agent produced all math,
 257 algorithms, and appendices.

258 Reproducibility Statement

259 We provide: (1) a deterministic theory-only reference implementation (typed executors, certificate
 260 calculators, calibration scripts) with seeds and container specs; (2) a governance table of sources,
 261 versions, cut dates for catalogue features; (3) a CI script that re-computes all certificate values for
 262 the micro-example and verifies *zero* leakage events. No external web access or wet-lab results are
 263 required to reproduce the paper’s numbers.

264 A Proofs

265 **Proof of Theorem 1.** By multiplicativity of the gated $\widehat{\text{PathConf}}$, invariance of $\widehat{\text{Rev}}, \widehat{\text{TFI}}$ to single-
 266 edge removal, and non-circular recomputation with $w_3 = 0$, one obtains the bound with C_P as
 267 defined.

268 **Proof of Theorem 2.** Under typed d -sep, shortcut evidence cannot appear along admissible meta-
 269 paths. Choose calibrated w^* and isotonic ϕ^* so that the certified terms dominate any residual
 270 $(\widehat{\text{Rev}}, \widehat{\text{TFI}})$ from spurious paths. Monotonicity preserves the strict margin.

271 **Proof of Theorem 3.** Let $g(P) = \widehat{\text{PathConf}}(P)$ be multiplicative over edges; removing edges
 272 multiplies g by factors ≤ 1 . The marginal decrement from removing e given S is smaller for larger S ,
 273 establishing diminishing returns and submodularity. Nonnegativity and normalization follow from
 274 construction.

275 B Micro-Example: Full Arithmetic

276 Augment Sec. 6 with a decoy path $P_{\text{spur}} : d \rightarrow x' \rightarrow y$ (missing Pathway). The typed executor
 277 rejects it; thus $\widehat{\text{PathConf}}(P_{\text{spur}}) = 0$. For the valid P , ablations over $\{d \rightarrow x, x \rightarrow p, p \rightarrow y\}$
 278 give normalized drops; the maximum is $\Delta\text{Score}_{\text{cf}}$. Using $w = (0.35, 0.2, 0.25, 0.2)$, $\widehat{\text{TFI}} = 0.5$,
 279 $\widehat{\text{Rev}} = 0.6$, $\text{IV} = 0.55$, and $\alpha = \beta = \gamma = \frac{1}{3}$, we obtain $\widehat{\text{MCS}} = \frac{1}{3}\Delta\text{Score}_{\text{cf}} + 0.35$. If the largest
 280 ablation is $x \rightarrow p$ (schema-violating), then $s' = 0.2 \cdot 0.6 + 0.2 \cdot 0.5 = 0.22$; with $\widehat{\text{PathConf}} = 0.504$,
 281 the base $s_0 = 0.3964$ and $\Delta\text{Score}_{\text{cf}} = 0.1764$.

282 C Governance Table Template

283 A machine-readable (CSV/JSON) table listing: source name, version, cut date t_0 , parser checksum,
 284 schema checksum, certificate script checksum, and CI status ($0/\langle N \rangle$ leakage violations). The CI
 285 recipe replays the micro-example, regenerates all values deterministically, and verifies checksums.