

# Transition-Matrix Regularization for Next Dialogue Act Prediction in Counselling Conversations

Anonymous ACL submission

## Abstract

This paper studies how empirical dialogue-flow statistics can be incorporated into Next Dialogue Act Prediction (NDAP). A KL regularization term is proposed that aligns predicted act distributions with corpus-derived transition patterns. Evaluated on a 60-class German counselling taxonomy using 5-fold cross-validation, this improves macro-F1 by 9–42% relative depending on encoder and substantially improves dialogue-flow alignment. Cross-dataset validation on HOPE suggests that improvements transfer across languages and counselling domains. In systematic ablations across pre-trained encoders and architectures, the findings indicate that transition regularization provides consistent gains and disproportionately benefits weaker baseline models. The results suggest that lightweight discourse-flow priors complement pre-trained encoders, especially in fine-grained, data-sparse dialogue tasks.

## 1 Introduction

NDAP forecasts the communicative function of the *upcoming* utterance from the dialogue history. Although the task has a long tradition in dialogue research (Nagata and Morimoto, 1994; Stolcke et al., 2000; Reithinger et al., 1996), it has received less attention in the era of large language models (LLMs). Yet it offers a structured and interpretable mechanism for steering LLM behavior, complementing prompting (Brown et al., 2020), instruction tuning (Wei et al., 2022), and reward-based approaches (Ouyang et al., 2022). By anticipating the next dialogue action, systems can condition prompts or constraints to encourage more stable, coherent, and goal-directed behavior (Chen et al., 2023).

In counselling and other high-structure domains, next acts follow consistent pragmatic patterns: greetings typically precede problem statements, exploration precedes intervention, and closing behaviors follow resolution (Bickmore et al., 2013; Al-

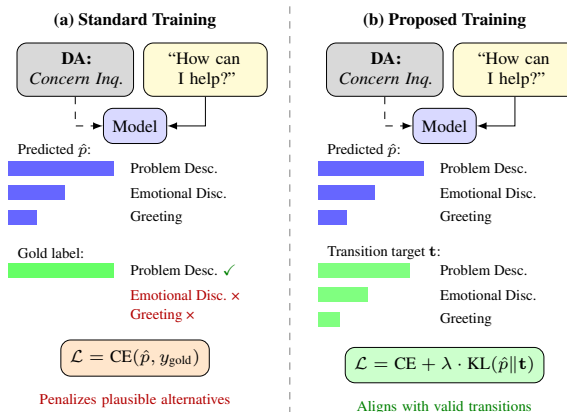


Figure 1: Comparison of standard NDAP training (a) vs. transition-matrix regularized training (b). Standard cross-entropy penalizes all non-gold predictions equally, even when multiple next acts are plausible. The proposed regularizer aligns predictions with empirical transition patterns from the corpus.

thoff et al., 2016). Classical dialogue managers explicitly modeled these transitions through Markov or CRF-based structures (Stolcke et al., 2000; Zimmermann, 2009). Modern neural systems, however, largely abandon symbolic transition models and instead rely on end-to-end architectures to infer discourse structure implicitly (Ultes et al., 2017; Ravuru et al., 2022).

This shift removes an inductive bias. Neural models see only a single gold next-act label per instance, providing limited signal when several next acts are plausible—a common situation in counselling (Wu et al., 2022b; Demasi et al., 2020). The gold label in NDAP is inherently under-specified: it represents one observed continuation among many valid possibilities. Standard cross-entropy supervision thus penalizes the model for predicting other plausible acts. Consequently, models may struggle to capture the multi-modal distribution of plausible next actions (Zhao et al., 2017).

This limitation is addressed by incorporating an

empirical transition matrix directly into the loss. The **transition-matrix KL regularizer** encourages the predicted next-act distribution to align with observed transition statistics, injecting pragmatic discourse-flow information as a soft, differentiable constraint (Figure 1). This preserves the flexibility of neural encoders while reinstating a structural prior reminiscent of classical systems.

This idea is evaluated in German text-based counselling, where communicative actions are fine-grained and governed by psychosocial norms. The dataset uses a five-level taxonomy with 60 dialogue act categories (Albrecht et al., 2025). NDAP is performed across all speaker transitions. To exploit the taxonomy structure, category history augmented architectures are introduced.

The results show that transition-based regularization provides consistent gains and disproportionately benefits weaker models. The regularizer is lightweight, model- and architecture-agnostic, and can be integrated as a drop-in objective without architectural modifications or sequence-level decoding. Improvements span both predictive metrics (F1, Top-3) and structural measures of dialogue-flow alignment.

These results show that domain-specific structural priors can enhance NDAP and offer a mechanism for steering LLM-driven systems in specialized domains (see Appendix A.2.1 for client-side-specific evaluation relevant to client simulation). The annotated counselling dataset is released to support further research on controllable and structured dialogue modelling.

## 2 Related Work

Dialogue act (DA) prediction is a well-established task that assigns communicative functions to observed utterances. Classical systems model discourse structure using stochastic grammars, HMMs, or CRFs (Stolcke et al., 2000; Geertzen, 2009), while recent neural approaches employ hierarchical encoders, contextual attention, and multi-modal cues (Colombo et al., 2020). These models capture local and long-range structure but operate entirely on *observed* utterances.

In contrast, our work addresses the fundamentally harder task of NDAP: forecasting the communicative function of the *upcoming* utterance without access to its surface form. Early statistical work showed that DA sequences exhibit strong structural regularities, with n-gram models improving

prediction through explicit transition constraints (Reithinger et al., 1996; Geertzen, 2009). Neural NDAP research remains limited. Prior work integrates multi-turn history using hierarchical attention (Tanaka et al., 2019), incorporates multimodal features, or employs semi-supervised consistency objectives (He et al., 2022). Latent-variable architectures (Ji et al., 2016) regularize discourse trajectories, yet rely on implicitly learned transitions. Recent studies in specialized domains (e.g., talk-move prediction, counselling) forecast future strategies or acts as auxiliary signals (Ganesh et al., 2021; Wu et al., 2022b), but do not incorporate explicit transition-based priors. Overall, existing NDAP approaches condition on DA history but do not constrain predictions using empirical dialogue-flow statistics.

Explicit modeling of DA transitions is well explored in *sequence labeling*, where the goal is to assign a DA label to every utterance in a conversation. Classical systems encode transitions through n-gram dialogue grammars or HMMs (Stolcke et al., 2000), and neural architectures commonly add CRF layers to impose label-transition structure (Chen et al., 2018; Shang et al., 2020). These methods learn transition potentials or decode full DA sequences with structural constraints. However, they do not address NDAP, where only the *next* act must be predicted and no sequence decoder is used. Critically, none incorporate an *empirical* transition matrix directly into the loss for distribution-level regularization.

Distribution-level constraints have been introduced through posterior regularization, which biases models toward constraint-satisfying priors using KL divergence (Ganchev et al., 2010). Such techniques have improved dialogue understanding and state tracking (Jin et al., 2018), and future-aware constraints have been applied to generation models (Feng et al., 2020). These works demonstrate the utility of KL-based structural guidance, but they do not incorporate DA-transition statistics nor target NDAP.

Recent approaches increasingly rely on LLMs for dialogue generation, typically steered via prompting or instruction tuning. However, growing evidence suggests that LLMs do not reliably acquire human-like discourse behavior, particularly with respect to pragmatic sequencing and role consistency (Shukuri et al., 2023; Wagner and Ultes, 2024; Rudolph et al., 2025). Fluent surface realization therefore provides limited leverage for

structured dialogue control, motivating the use of explicit structural priors.

Our approach fills a gap between these research threads. Unlike NDAP models that rely on implicitly learned transitions, we introduce an explicit, data-derived *transition-matrix regularizer* that aligns the predicted next-act distribution with empirical dialogue-flow patterns. Unlike CRF or HMM models, our method does not require sequence decoding. And unlike posterior-regularization approaches, our structural prior is grounded directly in observed DA transitions. To our knowledge, this is the first instance of integrating empirical DA-transition constraints into the optimization objective for NDAP.

### 3 Method

#### 3.1 Problem Formulation

Given a conversation history consisting of  $n$  utterances  $\{u_1, u_2, \dots, u_n\}$ , the goal is to predict the category  $c_t$  of the next utterance. All speaker transitions are considered (Counselor→Counselor, Counselor→Client, Client→Client, Client→Counselor), performing NDAP based on the conversation history. Categories belong to a five-level hierarchy with 60 leaf-level dialogue acts.

The task assumes access to gold DA labels in conversation history, appropriate for post-hoc analysis, training simulations, and constrained LLM generation.

To ground the task in a concrete taxonomy, the full hierarchical structure is adopted from the OnCoCo dataset, which organizes 60 leaf-level categories into progressively more abstract semantic groups (Albrecht et al., 2025). The hierarchy captures both conversational function and pragmatic counseling flow: high-level groups distinguish phases such as greetings, problem exploration, motivation building, and closing, while lower levels specify fine-grained communicative functions such as factual problem descriptions, emotional disclosures, resource activation, or evaluation of solution attempts.

OnCoCo provides finer granularity (60 categories vs. 3–15 in alternatives such as Anno-MI (Wu et al., 2022a), HOPE (Malhotra et al., 2022), or MITI (Moyers et al., 2016)) and integrates multiple counselling paradigms rather than a single therapeutic approach, critical for controllable client simulation.

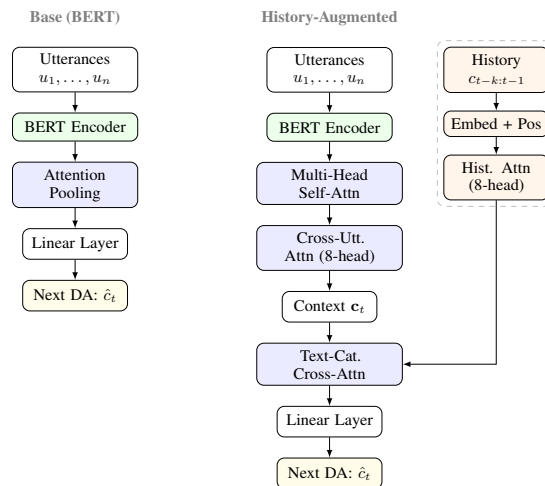


Figure 2: Model architectures for NDAP: BERT + attention pooling (left) and a history-augmented variant that incorporates conversational context and previous dialogue-act labels (right).

#### 3.2 Architectures

The base model encodes utterances via a pretrained BERT encoder and aggregates token representations through learned attention pooling. The pooled representation is passed to a linear classifier for NDAP.

The enhanced variant replaces attention pooling with multi-head self-attention for utterance encoding and adds 8-head cross-utterance attention to model conversation flow. A sliding window of the previous  $h$  categories is embedded with positional encodings and processed through 8-head self-attention. A text-category cross-attention mechanism integrates the historical context with the conversation representation, and a 4-head category transition attention layer models dependencies between consecutive acts (Figure 2).

#### 3.3 Training

Models combine cross-entropy loss with transition matrix loss regularization:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}} + \lambda_{\text{tm}} \mathcal{L}_{\text{TM}}. \quad (1)$$

Cross-entropy provides the primary supervision signal based on the single annotated next category, while the transition-matrix loss introduces a soft prior reflecting the empirical multi-modal distribution of plausible next acts.

#### 3.4 Transition Matrix and Dialogue Flow Regularization

A central component in our approach is the *transition matrix loss*, which encourages predictions

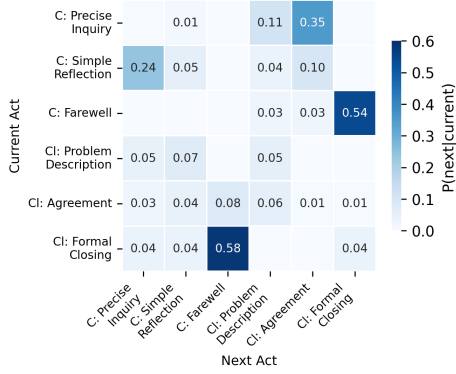


Figure 3: Example 6x6 subset (Fold 0) of the empirical transition matrix. C = Counselor, Cl = Client. High probabilities show pragmatic patterns: closings trigger closings, precise inquiry lead to agreement.

to respect observed category transition patterns in conversation. Although prior dialogue-act models have encoded transitions implicitly through sequential architectures or explicitly through statistical models, to the authors’ knowledge no prior work incorporates an empirical transition matrix directly into the training objective of a neural NDAP model.

### 3.4.1 Computing the Empirical Transition Matrix

A transition matrix encodes the probability distribution over next states given the current state—here, the likelihood of each dialogue act following another. From the training corpus, we compute a normalized empirical transition matrix  $\mathbf{T}$  where  $T_{ij}$  represents the probability of transitioning from category  $i$  to category  $j$ . Sparse transitions are smoothed with additive smoothing to avoid zero probabilities. In social counselling conversations, the transition matrix exhibits clear pragmatic patterns as can be seen in figure 3.

### 3.4.2 Transition Matrix Loss

The transition matrix loss uses KL divergence to measure alignment between model predictions and empirical category transitions. Given the previous category  $c_{t-1}$ , the target transition distribution is:

$$\mathbf{t}^{(t)} = \mathbf{T}[c_{t-1}] \in \mathbb{R}^C \quad (2)$$

The model’s predicted distribution is:

$$\hat{\mathbf{p}}^{(t)} = \text{softmax}(\hat{\mathbf{y}}_{\text{final}}) \quad (3)$$

The transition matrix loss measures divergence:

$$\mathcal{L}_{\text{TM}} = \text{KL}(\hat{\mathbf{p}}^{(t)} \parallel \mathbf{t}^{(t)}) = \sum_{j=1}^C \hat{p}_j^{(t)} \log \frac{\hat{p}_j^{(t)}}{t_j^{(t)}} \quad (4)$$

KL divergence is asymmetric, penalizing impossible transitions more heavily than missing low-probability valid ones, making it well-suited for enforcing dialogue structure.

Our transition-matrix regularizer differs from label smoothing (Szegedy et al., 2015), which distributes probability mass uniformly across non-target classes. In contrast, our approach distributes mass according to empirically observed transition patterns, providing domain-specific rather than uniform smoothing (Table 4; see Appendix A.2.2 for additional comparisons on the client simulation subset).

The weight  $\lambda_{\text{tm}} \in \{0.0, 0.2, 0.5, 1.0, 1.5\}$  is explored systematically (Section 4.4).

## 4 Experiments

### 4.1 Dataset

Models are evaluated on a corpus of German social counselling dialogues consisting of 76 conversations with 5,457 utterances. The data originate from structured counselling role-play sessions conducted by social science students in a university course on online text-based counselling. Participants acted in predefined client and counsellor roles, and no real clients or personal information were involved. All conversations were anonymized and collected with consent for research use.

Each utterance is annotated with a category from the OnCoCo taxonomy (Albrecht et al., 2025), which defines 60 leaf-level dialogue act categories organized in a five-level hierarchy. While the taxonomy has been previously described, this work contributes the first publicly available corpus of complete counselling conversations annotated with this scheme, enabling sequential modeling of dialogue flow. Annotations were produced by trained social science students and subsequently reviewed by a domain expert. This sequential review process precludes formal inter-rater agreement but ensures annotation quality through expert oversight.

The human-annotated conversational dataset, along with metadata describing the annotation schema, is released to support reproducibility and further research on hierarchical dialogue-act pre-

diction and counselling dialogue modeling.<sup>1</sup>

For experimentation, NDAP is performed across all speaker transitions, yielding instances across 60 categories. Evaluation uses 5-fold cross-validation with conversation-level partitioning to prevent data leakage. For each fold, training is conducted on 80% of conversations with evaluation on the held-out 20%. Mean performance  $\pm$  standard deviation across all five test folds is reported. Rather than selecting optimal hyperparameters on a separate validation set (which would further reduce already limited training data), results across the full hyperparameter grid are reported (Section 4.4), allowing readers to assess performance at each regularization strength. The original class distribution exhibits substantial imbalance. LLM-based synthetic data augmentation was also explored to address this imbalance; however, it did not improve BERT-based models (see Appendix A.2.3 for details).

## 4.2 Baselines and Models

The history-augmented architecture is compared against several baseline approaches: (1) a transition-matrix baseline that uses only empirical dialogue-flow patterns without learning, (2) a Simple RNN baseline, and (3) the architecture proposed by Tanaka et al. (2019), which uses hierarchical attention to integrate multi-turn dialogue history for NDAP.

To validate the robustness of the findings across different pretrained language models, all neural baselines and history-aware models were tested using 7 different German BERT variants: EuroBERT-210m, EuroBERT-610m, G-BERT-base, G-BERT-large, GELECTRa-base, Modern-G-BERT-134M, and Modern-G-BERT-1B. Additionally, context window size is varied (1, 4, 8, 12 utterances) and transition-matrix regularization is compared against label smoothing ( $\epsilon \in \{0.0, 0.1, 0.2\}$ ) as an alternative regularization strategy (see Appendix A.2.2). This systematic evaluation across 300 configurations (7 encoders  $\times$  2 architectures  $\times$  5 TM weights  $\times$  4 context lengths, plus label smoothing variants), evaluated using 5-fold cross-validation, shows that architectural improvements hold consistently across language models ranging from 110M to 1B parameters. Table 1 summarizes the model architectures and their properties.

**Transition Matrix Baseline:** The transition matrix baseline provides a competitive non-learning

| Model              | Attention   | History |
|--------------------|-------------|---------|
| Transition Matrix  | —           | —       |
| GPT-5-mini (LLM)   | —           | ✓       |
| gpt-oss-120b (LLM) | —           | ✓       |
| Simple RNN         | Pooling     | —       |
| Tanaka (2019)      | Hier. Attn. | ✓       |
| BERT               | Pooling     | —       |
| BERT+History       | Multi-head  | ✓       |

Table 1: Model architectures and properties.

baseline. This model predicts the most likely next category given only the previous utterance’s category, using the empirical transition matrix computed from training data. It requires no neural training or context encoding and serves as a practical reference point for the value of learned contextual representations.

**RNN Baselines:** Two RNN-based baselines are implemented: (1) Simple RNN, a basic recurrent architecture over utterance embeddings, and (2) the Tanaka et al. (2019) architecture, which uses hierarchical attention mechanisms to model multi-turn context for NDAP. All RNN baselines are trained with the same transition-matrix regularization as BERT models. Results are shown in Section 4.4.1.

**LLM Baseline:** To contextualize fine-tuned BERT performance against state-of-the-art language models, both proprietary (GPT-5-mini) and open-source (gpt-oss-120b, a 120B-parameter Mixture of Experts model) LLMs are evaluated in a zero-shot setting. Each model receives conversation history (last 12 turns) and all 60 category descriptions, returning top-3 predictions with confidence scores. This baseline tests whether large-scale pretraining alone captures dialogue-flow patterns without explicit structural constraints. The full prompt template is provided in Appendix A.2.7.

## 4.3 Evaluation Metrics

Models are evaluated on predictive correctness and dialogue-flow alignment via three metric categories:

**Predictive Correctness** Macro-F1, weighted F1, and Top-3 accuracy serve as the primary measures. Top-3 accuracy acknowledges that multiple next acts can be plausible. Note that these metrics assume a single correct answer, which is a simplification: in NDAP, multiple dialogue acts are often genuinely valid continuations, so moderate absolute scores are expected.

<sup>1</sup>Code and data available at <https://anonymous.4open.science/r/tm-reg-8E74>

| Encoder          | Best $_{\lambda=0}$ | Best $_{\lambda>0}$ | $\lambda$ | $\Delta$ | %Gain  |
|------------------|---------------------|---------------------|-----------|----------|--------|
| GBERT-large      | .097                | <b>.102</b>         | 0.5       | +0.005   | +5.1%  |
| GBERT-base       | .092                | .098                | 0.5       | +0.005   | +6.0%  |
| ModernGBERT-1B   | .089                | .096                | 0.2       | +0.007   | +8.5%  |
| EuroBERT-610M    | .087                | .096                | 0.5       | +0.009   | +10.5% |
| EuroBERT-210M    | .080                | .095                | 0.5       | +0.014   | +18.0% |
| ModernGBERT-134M | .076                | .086                | 0.5       | +0.009   | +12.2% |
| GELECTRa-base    | .060                | .070                | 1.0       | +0.009   | +16.4% |
| <i>Mean</i>      | .083                | .092                | -         | +0.008   | +11.0% |

Table 2: TM regularization effect by encoder (macro-F1, 60 categories, 5-fold CV). Compares best result at  $\lambda_{tm}=0$  vs best at  $\lambda_{tm}>0$ . Sorted by best macro-F1.

**Dialogue-Flow Alignment** To assess how well model predictions adhere to the conversational structure observed in the training data, a group of metrics based on a first-order empirical transition matrix  $T$  is employed: Cumulative Accuracy at 70% (Cum70) and Jensen-Shannon (JS) Divergence. Cum70 measures whether the predicted category falls within the set of most likely transitions that together account for 70% of the empirical probability mass—capturing whether predictions align with pragmatically plausible continuations. JS divergence measures the overall distributional alignment between predicted and empirical transition distributions.

**Rationale for First-Order Metrics** While dialogue context is inherently rich and multi-turn, the dialogue-flow metrics deliberately rely on a first-order (last-act-to-next-act) transition matrix. This is a practical necessity. Constructing higher-order transition matrices (e.g., second-order, based on the last two acts) would be infeasible given the 60 leaf categories, as it would lead to a combinatorial explosion of states ( $60^2 = 3600$  potential source pairs) and result in an extremely sparse and unreliable matrix with this dataset size. The first-order matrix thus serves as a robust and computable proxy for measuring the model’s grasp of foundational, short-term dialogue coherence.

**Transition Matrix Computation** The transition matrix is computed from training data within each CV fold, measuring whether the model internalized valid dialogue flows.

## 4.4 Results

### 4.4.1 TM Regularization Effect by Encoder

Table 2 shows the effect of transition-matrix regularization on the 60-category classification task,

| $\lambda_{tm}$ | Macro-F1               | W-F1                   | Top-3                  | Cum70                  | JS                     |
|----------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| 0.0            | .309 $\pm$ .011        | .495 $\pm$ .008        | <b>.789</b> $\pm$ .011 | .913 $\pm$ .019        | .299 $\pm$ .017        |
| 0.2            | .314 $\pm$ .013        | .496 $\pm$ .010        | .783 $\pm$ .014        | .915 $\pm$ .016        | .258 $\pm$ .013        |
| 0.5            | <b>.319</b> $\pm$ .014 | .499 $\pm$ .012        | .781 $\pm$ .013        | .916 $\pm$ .014        | .225 $\pm$ .011        |
| 1.0            | .315 $\pm$ .015        | .499 $\pm$ .012        | .777 $\pm$ .013        | <b>.918</b> $\pm$ .013 | .205 $\pm$ .008        |
| 1.5            | .317 $\pm$ .013        | <b>.501</b> $\pm$ .010 | .773 $\pm$ .011        | .916 $\pm$ .009        | <b>.201</b> $\pm$ .007 |

Table 3: Cross-dataset validation on HOPE (15 English counselling dialogue act classes, 5-fold CV, mean  $\pm$  std).

comparing each encoder’s best result at  $\lambda_{tm}=0$  versus its best result at  $\lambda_{tm}>0$ .

All encoders improve with transition-loss regularization (Table 2). Relative gains range from +5.1% (GBERT-large) to +18.0% (EuroBERT-210M), with a mean improvement of +11.0%. Notably, smaller encoders show larger relative gains, suggesting TM regularization partially compensates for limited model capacity. Optimal  $\lambda_{tm}$  values cluster around 0.5, with only GELECTRa-base benefiting from higher regularization ( $\lambda_{tm}=1.0$ ). Architecture effects vary by encoder (Appendix B Table 16). Figure 4 visualizes values averaged across architectures and context sizes for each  $\lambda_{tm}$  value; this averaging reveals larger relative improvements, ranging from 9% (GBERT-large) to 42% (GELECTRa-base), consistent with the abstract.

### 4.4.2 Cross-Dataset Validation

To validate generalization beyond German counselling, the transition-matrix regularizer is evaluated on the HOPE dataset (Malhotra et al., 2022), an English counselling corpus with 15 dialogue act categories. Table 3 summarizes results using XLM-RoBERTa-base and BERT-base encoders with both BERT and History architectures. Additional evaluation on Switchboard (SWDA), a non-counselling benchmark with highly skewed transition distributions, is provided in Section A.2.4.

Results indicate generalization across languages (German to English), counselling modalities (online text-based to spoken), and category systems (60-class OnCoCo taxonomy to 15-class HOPE scheme). Macro-F1 improves from 0.309 to 0.319 (+3.2% relative) at  $\lambda_{tm}=0.5$ , and JS divergence drops from 0.299 to 0.201 (33% reduction).

### 4.4.3 Label Smoothing Comparison

Table 4 compares transition-matrix regularization against label smoothing (Szegedy et al., 2015) across all seven encoders. TM regularization con-

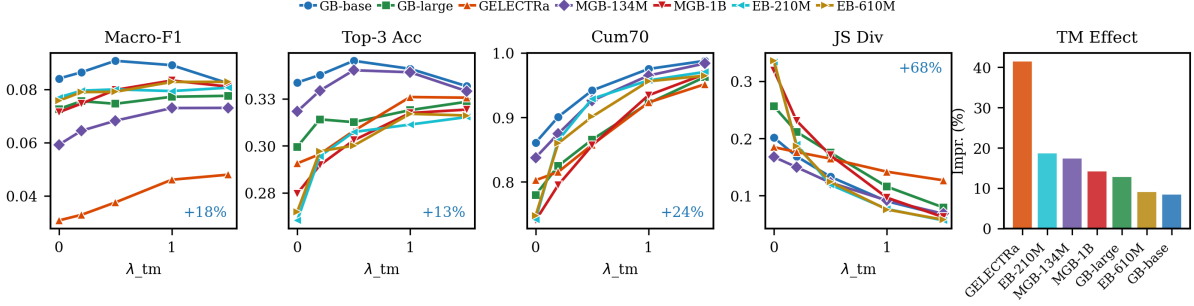


Figure 4: Effect of transition loss weight ( $\lambda_{tm}$ ) on 60-category classification. Left four panels show how Macro-F1, Top-3 Accuracy, Cum70, and JS Divergence change with increasing  $\lambda_{tm}$  across encoders. Right panel compares relative macro-F1 gains from TM regularization by encoder.

| Encoder          | Label Smoothing |                |                | TM Reg.     |          |
|------------------|-----------------|----------------|----------------|-------------|----------|
|                  | $\epsilon=0.0$  | $\epsilon=0.1$ | $\epsilon=0.2$ | Best        | $\Delta$ |
| GBERT-large      | .079            | .074           | .069           | <b>.102</b> | +0.023   |
| GBERT-base       | .073            | .081           | .071           | <b>.098</b> | +0.017   |
| ModernGBERT-1B   | .059            | .074           | .065           | <b>.096</b> | +0.022   |
| EuroBERT-610M    | .085            | .078           | .077           | <b>.096</b> | +0.010   |
| EuroBERT-210M    | .075            | .066           | .067           | <b>.095</b> | +0.019   |
| ModernGBERT-134M | .049            | .050           | .047           | <b>.086</b> | +0.035   |
| GELECTRa-base    | .033            | .033           | .031           | <b>.070</b> | +0.036   |
| <i>Mean</i>      | .065            | .065           | .061           | <b>.092</b> | +0.023   |

Table 4: Label smoothing vs. TM regularization (macro-F1, 60 categories, 5-fold CV).  $\Delta$  = TM best – best LS.

484 sistently outperforms label smoothing on all encoders, with a mean improvement of +0.023 macro-  
 485 F1 over the best label smoothing configuration.  
 486 The relative gains are largest for weaker models  
 487 (GELECTRa-base: +0.036, ModernGBERT-134M:  
 488 +0.035), confirming that transition-based priors pro-  
 489 vide the strongest benefits when baseline perfor-  
 490 mance is low.  
 491

#### 4.4.4 Model Comparison and Effect of Regularization Strength

492 Figure 4 visualizes the effect of transition loss  
 493 weight ( $\lambda_{tm}$ ) across architectures and metrics.  
 494 Dialogue-flow metrics (Cum70, JS divergence) im-  
 495 prove monotonically with increasing  $\lambda_{tm}$ , while  
 496 predictive metrics (Macro-F1, Top-3) peak around  
 497  $\lambda_{tm}=1.0$ – $1.5$ . The right panel shows that TM regu-  
 498 larization provides larger effect sizes than architec-  
 499 ture changes (BERT→History).  
 500

501 Table 5 compares performance across seven Ger-  
 502 man BERT variants (110M–1B parameters), two  
 503 LLMs, and baselines. A transition matrix base-  
 504 line that predicts using only empirical  $P(\text{next}|\text{prev})$   
 505 from training data achieves macro-F1 of 0.056,  
 506 outperforming RNN models (0.003–0.008) but  
 507

| Model   | $\lambda_{tm}$ | Cfg | Macro-F1            | W-F1                | Top-3               | Cum70               | JS                  |
|---|----------------|-----|---------------------|---------------------|---------------------|---------------------|---------------------|
| <i>LLM Baselines (zero-shot, 5-fold CV)</i>         |                |     |                     |                     |                     |                     |                     |
| GPT-5-mini  | -              | -   | .091 ± 0.011        | .132 ± 0.014        | .254 ± 0.014        | .550 ± 0.035        | -                   |
| gpt-oss-120b  | -              | -   | .072 ± 0.006        | .108 ± 0.007        | .174 ± 0.010        | .400 ± 0.021        | -                   |
| <i>RNN Baselines (5-fold CV)</i>                    |                |     |                     |                     |                     |                     |                     |
| Simple RNN  | 0.5            | -   | .003 ± 0.000        | .014 ± 0.002        | .208 ± 0.012        | .751 ± 0.037        | .198 ± 0.006        |
| Tanaka (2019)                                       | 1.5            | -   | .008 ± 0.002        | .031 ± 0.008        | .235 ± 0.017        | .802 ± 0.035        | .209 ± 0.015        |
| <i>Transition Matrix Baseline (5-fold CV)</i>       |                |     |                     |                     |                     |                     |                     |
| TM Only   | -              | -   | .056 ± 0.003        | .115 ± 0.011        | .315                | 1.00                | -                   |
| <i>Fine-tuned Encoders (best config, 5-fold CV)</i> |                |     |                     |                     |                     |                     |                     |
| GB-large  | 0.5            | H4  | <b>.102 ± 0.007</b> | <b>.171 ± 0.012</b> | .342 ± 0.023        | .894 ± 0.013        | .175 ± 0.004        |
| GB-base   | 0.5            | B8  | .098 ± 0.010        | .158 ± 0.015        | <b>.354 ± 0.024</b> | .933 ± 0.008        | .143 ± 0.006        |
| MGB-1B  | 0.2            | H8  | .097 ± 0.011        | .159 ± 0.021        | .314 ± 0.026        | .821 ± 0.022        | .235 ± 0.015        |
| EB-610M   | 0.5            | H4  | .097 ± 0.012        | .160 ± 0.019        | .334 ± 0.024        | .939 ± 0.007        | .109 ± 0.007        |
| EB-210M   | 0.5            | H12 | .095 ± 0.015        | .158 ± 0.020        | .331 ± 0.027        | .937 ± 0.011        | .125 ± 0.004        |
| MGB-134M  | 0.5            | H4  | .086 ± 0.007        | .147 ± 0.009        | .358 ± 0.021        | <b>.942 ± 0.005</b> | .122 ± 0.001        |
| GELECTRa  | 1.0            | H4  | .070 ± 0.007        | .128 ± 0.011        | .335 ± 0.015        | .933 ± 0.019        | <b>.121 ± 0.003</b> |

Table 5: Results by model (60 categories, 5-fold CV). Cfg = best architecture (B=BERT, H=History) + context length.

508 falling well short of fine-tuned encoders (0.070–  
 509 0.102). Among fine-tuned encoders, GBERT-  
 510 large achieves the highest macro-F1 (0.102), while  
 511 smaller models like ModernGBERT-134M achieve  
 512 competitive performance with TM regularization.  
 513 The Cfg column indicates the best architecture  
 514 (B=BERT, H=History) and context length for each  
 515 encoder. History-augmented models dominate,  
 516 with optimal context lengths of 4–12 utterances.  
 517 LLMs (GPT-5-mini, gpt-oss-120b) achieve lower  
 518 dialogue-flow alignment (Cum70: 0.40–0.55) de-  
 519 spite competitive macro-F1, indicating that pre-  
 520 training alone does not capture transition patterns.

521 **Significance Testing.** Paired bootstrap tests with  
 522 Benjamini-Hochberg correction confirm that TM  
 523 regularization yields robust but configuration-  
 524 dependent improvements, with mid-sized encoders  
 525 showing the most consistent gains (Table 14).  
 526 History-based architectures provide benefits only  
 527 for specific encoders and do not consistently outper-  
 528 form standard BERT (see Appendix B for encoder-  
 529 specific results).

## 5 Discussion

The results highlight several insights:

**1. Transition regularization provides consistent benefits:** The transition-matrix regularizer improves all metrics across configurations (Table 4). This dual improvement in predictive accuracy and dialogue-flow alignment reflects the multi-modal nature of NDAP: cross-entropy treats all non-gold predictions as equally wrong, while TM regularization provides *distributional* supervision encoding plausible next acts. Moreover, the proposed regularizer is genuinely model- and architecture-agnostic: it yields consistent gains across all seven encoder variants (110M–1B parameters) and both attention-pooling and history-augmented architectures, and can be integrated as a drop-in objective without any architectural changes or sequence-level decoding.

**2. Structural priors matter more than model scale:** Across a  $10\times$  size range (110M–1B parameters), TM regularization provides larger gains than encoder choice (Table 5). ModernGBERT-1B achieves best macro-F1 at  $\lambda_{tm}=0.2$  but exhibits lower dialogue-flow alignment (Cum70=.822, JS=.236) than smaller models with higher  $\lambda_{tm}$ —larger models may require less aggressive regularization but still benefit from structural priors.

**3. Differential benefits across encoders and architectures:** Encoder-level analysis (Appendix B) reveals a clear separation. TM regularization acts as a general-purpose inductive bias that disproportionately benefits weaker encoders (GELECTRa-base: +0.036, ModernGBERT-134M: +0.035 macro-F1), whereas explicit history modeling provides selective gains that do not generalize across encoder families (GELECTRa: +1.93 points; GBERT-base: -0.16 points). Across all encoders, domain-specific transition priors consistently outperform uniform label smoothing (Table 4).

**4. LLM comparison validates transition regularization:** GPT-5-mini achieves lower macro-F1 than the best fine-tuned model Table (5) and substantially lower dialogue-flow alignment (Cum70: 0.550 vs 0.894). The open-source gpt-oss-120b performs worse (macro-F1: 0.072) due to unreliable schema adherence in structured outputs, resulting in a 17% parse error rate. The low Cum70 for both LLMs confirms that even state-of-the-art LLMs do not naturally acquire dialogue-flow transition patterns from pretraining alone, validating the need for explicit structural priors.

**5. Client simulation subset shows strong benefits:** On a 28-category client-side subset (Appendix A.2.1), TM regularization yields +18% weighted F1 improvement (0.225→0.265) confirming that the approach is effective for the client simulation use case.

### 5.1 Effect of Regularization Strength

**Response pattern:** All metrics improve with transition-loss weight, with dialogue-flow metrics (cumulative coverage, JS divergence) continuing to improve at higher  $\lambda_{tm}$  values while predictive metrics peak around  $\lambda_{tm}=1.0$ –1.5. The increasing trend suggests the regularizer provides a stable optimization signal without creating competing objectives.

**Implicit regularization effect:** The transition-matrix regularizer also acts as an implicit regularizer against overfitting. Without it ( $\lambda_{tm}=0.0$ ), 83% of runs show overfitting magnitude  $>0.5$ . Positive  $\lambda_{tm}$  reduces this: at  $\lambda_{tm}=1.5$ , overfitting magnitude falls to 0.20 with no runs exceeding 0.5.

**Practical recommendation:** For practitioners, a  $\lambda_{tm} \in [0.5, 1.0]$  as a default is recommended. Across all configurations,  $\lambda_{tm}=0.5$  most frequently achieves optimal macro-F1, while  $\lambda_{tm}=1.0$  performs within 5% of optimal in over 70% of configurations, making it a robust choice when extensive tuning is not feasible. This pattern also holds for HOPE.

## 6 Conclusion

This paper proposed a transition-matrix KL regularizer for NDAP that incorporates empirical dialogue-flow structure into the training objective. Across experiments on a fine-grained German counselling taxonomy and cross-dataset transfer to other dialogue corpora, the regularizer consistently improves both predictive performance and alignment with observed dialogue-flow dynamics. Additionally, we presented history-augmented architectures that leverage broader dialogue context through multi-head attention.

A natural next step is integrating NDAP with LLM control, using predicted categories to condition prompts or guide sampling strategies for client simulation. Additionally, fine-tuning decoder-only architectures on NDAP could leverage their autoregressive nature to better model dialogue sequences while potentially benefiting from transition-matrix regularization.

## 629 Limitations

630 The primary dataset consists of 76 conversations  
631 with approximately 5,500 utterances—a relatively  
632 small corpus for training neural models. While  
633 cross-validation and cross-dataset experiments pro-  
634 vide evidence of robustness, the limited size re-  
635 stricts conclusions about generalization to larger-  
636 scale deployments or substantially different coun-  
637 selling contexts.

638 The conversations are role-play sessions con-  
639 ducted by social science students, not recordings  
640 of real clinical counselling. While participants fol-  
641 lowed structured scenarios designed to reflect au-  
642 thentic counselling dynamics, student role-plays  
643 may lack the pragmatic complexity, emotional  
644 depth, and unpredictability of genuine therapeu-  
645 tic interactions. Models trained on this data may  
646 not fully capture the nuances present in real-world  
647 counselling.

648 Annotations were produced by trained students  
649 and subsequently reviewed by a domain expert in  
650 a sequential process. This workflow precludes for-  
651 mal inter-rater agreement metrics, which are stan-  
652 dard for validating annotation reliability. Although  
653 expert review ensures quality control, the absence  
654 of quantified agreement limits assessment of anno-  
655 tation consistency and may affect reproducibility  
656 of the label assignments.

657 The transition matrix loss assumes transitions  
658 are stable across the corpus; in applications with  
659 significant domain shift or concept drift, alterna-  
660 tive regularization approaches may be needed. Our  
661 evaluation uses gold dialogue act labels for the con-  
662 versation history; although this is common practice  
663 in NDAP, these markings would have to be pre-  
664 dicted in fully autonomous operations, which could  
665 potentially lead to chain errors. End-to-end evalua-  
666 tion combining NDAP with LLM-based response  
667 generation remains unexplored.

## 668 Ethical considerations

669 The counselling dialogue data used in this study  
670 were collected from structured role-play sessions  
671 conducted by social science students in a univer-  
672 sity course setting. Participants acted in predefined  
673 client and counsellor roles; no real clients, patients,  
674 or personal therapeutic information were involved.  
675 All participants provided informed consent for re-  
676 search use of the data, and all conversations were  
677 anonymized prior to analysis. The dialogues were  
678 annotated by paid social science students with prior

679 training in counselling concepts. Annotation in-  
680 cluded both dialogue act labeling and anonymiza-  
681 tion of the conversations. All annotations were  
682 subsequently reviewed by a domain expert to en-  
683 sure consistency and quality.

684 We acknowledge that dialogue act prediction  
685 technology for counselling contexts could poten-  
686 tially be misused. However, the intended applica-  
687 tion, controllable client simulation for counselor  
688 training, serves an educational purpose that may  
689 improve the quality of counseling services. The  
690 data set released contains only simulated conversa-  
691 tions and poses minimal privacy risk.

## Acknowledgments

692 Generative AI tools (ChatGPT and Claude) were  
693 used throughout the article for coding assistance, to  
694 identify relevant literature, LaTeX figure and table  
695 preparation, and to improve clarity and style of  
696 writing. They were not used to generate scientific  
697 claims or experimental results. All content was  
698 reviewed, verified, and finalized by the authors.  
699

## References

- 700  
701 Jens Albrecht, Robert Lehmann, Aleksandra Polter-  
702 mann, Eric Rudolph, Philipp Steigerwald, and Mara  
703 Stieler. 2025. [OnCoCo 1.0: A Public Dataset  
704 for Fine-Grained Message Classification in On-  
705 line Counseling Conversations](#). *arXiv preprint*.  
706 ArXiv:2512.09804 [cs].
- 707 Tim Althoff, Kevin Clark, and Jure Leskovec. 2016. [708  
709 Large-scale Analysis of Counseling Conversations:  
710 An Application of Natural Language Processing to  
711 Mental Health](#). *Transactions of the Association for  
712 Computational Linguistics*, 4:463–476.
- 712 Timothy W. Bickmore, Daniel Schulman, and Candace  
713 Sidner. 2013. [Automated interventions for multiple  
714 health behaviors using conversational agents](#). *Patient  
715 Education and Counseling*, 92(2):142–148.
- 716 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
717 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
718 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
719 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
720 Gretchen Krueger, Tom Henighan, Rewon Child,  
721 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
722 Winter, and 12 others. 2020. [Language models are  
723 few-shot learners](#). In *Advances in neural information  
724 processing systems*, volume 33, pages 1877–1901.  
725 Curran Associates, Inc.
- 726 Sirui Chen, Yuan Wang, Zijing Wen, Zhiyu Li, Chang-  
727 shuo Zhang, Xiao Zhang, Quan Lin, Cheng Zhu,  
728 and Jun Xu. 2023. [Controllable Multi-Objective Re-  
729 ranking with Policy Hypernetworks](#). In *Proceedings*



|     |  |   |     |
|-----|--|---|-----|
| 842 | Eric Rudolph, Philipp Steigerwald, and Jens Albrecht.                        | Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin                             | 899 |
| 843 | 2025. <a href="#">Comparing human roleplayers and LLM-</a>                   | Guu, Adams Wei Yu, Brian Lester, Nan Du, An-                                  | 900 |
| 844 | <a href="#">simulated clients in online counselling training: An</a>         | drew M. Dai, and Quoc V. Le. 2022. <a href="#">Finetuned lan-</a>             | 901 |
| 845 | <a href="#">analysis of counselling patterns</a> . In <i>Proceedings of</i>  | <a href="#">guage models are zero-shot learners</a> . In <i>International</i> | 902 |
| 846 | <i>the 18th international conference on educational data</i>                 | <a href="#">conference on learning representations (ICLR) 2022</a> .          | 903 |
| 847 | <i>mining</i> , pages 381–387, Palermo, Italy. International                 |   |     |
| 848 | Educational Data Mining Society.   |   |     |
| 849 | Guokan Shang, Antoine Tixier, Michalis Vazirgiannis,                         | Zixiu Wu, Simone Balloccu, Vivek Kumar, Rim                                   | 904 |
| 850 | and Jean-Pierre Lorré. 2020. <a href="#">Speaker-change Aware</a>            | Helaoui, Ehud Reiter, Diego Reforgiato Recupero,                              | 905 |
| 851 | <a href="#">CRF for Dialogue Act Classification</a> . In <i>Proceedings</i>  | and Daniele Riboni. 2022a. <a href="#">Anno-MI: A Dataset</a>                 | 906 |
| 852 | <i>of the 28th International Conference on Computa-</i>                      | <a href="#">of Expert-Annotated Counselling Dialogues</a> . In                | 907 |
| 853 | <i>tional Linguistics</i> , pages 450–464, Barcelona, Spain                  | <i>ICASSP 2022 - 2022 IEEE International Confer-</i>                          | 908 |
| 854 | (Online). International Committee on Computational                           | <i>ence on Acoustics, Speech and Signal Processing</i>                        | 909 |
| 855 | Linguistics.   | <i>(ICASSP)</i> , pages 6177–6181, Singapore, Singapore.                      | 910 |
|     |  | IEEE.   | 911 |
| 856 | Elizabeth Shriberg, Rebecca Bates, Paul Taylor, An-                          | Zixiu Wu, Rim Helaoui, Diego Reforgiato Recupero,                             | 912 |
| 857 | dreas Stolcke, Daniel Jurafsky, Klaus Ries, Noah                             | and Daniele Riboni. 2022b. <a href="#">Towards Automated</a>                  | 913 |
| 858 | Coccaro, Rachel Martin, Marie Meteer, and Carol                              | <a href="#">Counselling Decision-Making: Remarks on Thera-</a>                | 914 |
| 859 | Van Ess-Dykema. 1998. Can prosody aid the auto-                              | <a href="#">pist Action Forecasting on the AnnoMI Dataset</a> . In            | 915 |
| 860 | matic classification of dialog acts in conversational                        | <i>Interspeech 2022</i> , pages 1906–1910. ISCA.                              | 916 |
| 861 | speech? <i>Language and Speech</i> , 41(3–4):439–487.                        |   |     |
| 862 | Kotaro Shukuri, Ryoma Ishigaki, Jundai Suzuki,                               | Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017.                          | 917 |
| 863 | Tsubasa Naganuma, Takuma Fujimoto, Daisuke                                   | <a href="#">Learning Discourse-level Diversity for Neural Dialog</a>          | 918 |
| 864 | Kawakubo, Masaki Shuzo, and Eisaku Maeda. 2023.                              | <a href="#">Models using Conditional Variational Autoencoders</a> .           | 919 |
| 865 | <a href="#">Meta-control of Dialogue Systems Using Large Lan-</a>            | In <i>Proceedings of the 55th Annual Meeting of the</i>                       | 920 |
| 866 | <a href="#">guage Models</a> . <i>arXiv preprint</i> . ArXiv:2312.13715      | <i>Association for Computational Linguistics (Volume</i>                      | 921 |
| 867 | [cs].  | <i>1: Long Papers)</i> , pages 654–664, Vancouver, Canada.                    | 922 |
|     |  | Association for Computational Linguistics.                                    | 923 |
| 868 | Andreas Stolcke, Klaus Ries, Noah Coccaro, Eliza-                            | Matthias Zimmermann. 2009. <a href="#">Joint segmentation and</a>             | 924 |
| 869 | beth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul                          | <a href="#">classification of dialog acts using conditional random</a>        | 925 |
| 870 | Taylor, Rachel Martin, Carol Van Ess-Dykema, and                             | <a href="#">fields</a> . In <i>Interspeech 2009</i> , pages 864–867. ISCA.    | 926 |
| 871 | Marie Meteer. 2000. <a href="#">Dialogue act modeling for au-</a>            |   |     |
| 872 | <a href="#">tomatic tagging and recognition of conversational</a>            | <b>A Model Architecture Details</b>   | 927 |
| 873 | <a href="#">speech</a> . <i>Computational Linguistics</i> , 26(3):339–374.   | <b>A.1 Hyperparameter Settings</b>  | 928 |
| 874 | Place: Cambridge, MA Publisher: MIT Press.                                   | <b>A.2 Ablation Study</b>   | 929 |
| 875 | Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe,                          | Ablation studies are conducted to validate the de-                            | 930 |
| 876 | Jonathon Shlens, and Zbigniew Wojna. 2015. <a href="#">Re-</a>               | sign choices: (1) evaluating on a client simulation                           | 931 |
| 877 | <a href="#">thinking the Inception Architecture for Computer</a>             | subset, (2) comparing transition-matrix regulariza-                           | 932 |
| 878 | <a href="#">Vision</a> . <i>arXiv preprint</i> . ArXiv:1512.00567 [cs].      | tion against label smoothing, (3) evaluating the                              | 933 |
| 879 | Koji Tanaka, Junya Takayama, and Yuki Arase. 2019.                           | effect of synthetic data augmentation, (4) exploring                          | 934 |
| 880 | Dialogue-act prediction of future responses based                            | alternative KL formulations, and (5) investigating                            | 935 |
| 881 | on conversation history. In <i>Proceedings of the 57th</i>                   | hierarchical label embeddings.  | 936 |
| 882 | <i>annual meeting of the association for computational</i>                   | <b>A.2.1 Client Simulation Subset</b>   | 937 |
| 883 | <i>linguistics: Student research workshop</i> , pages 197–                   | For client simulation applications, predicting                                | 938 |
| 884 | 202.   | client-side dialogue acts is the primary goal.                                | 939 |
| 885 | Stefan Ultes, Lina M. Rojas-Barahona, Pei-Hao Su,                            | TM regularization is evaluated on a subset con-                               | 940 |
| 886 | David Vandyke, Dongho Kim, Iñigo Casanueva,                                  | taining only client-side categories, filtering for                            | 941 |
| 887 | Paweł Budzianowski, Nikola Mrkšić, Tsung-Hsien                               | counselor→client and client→client transitions.                               | 942 |
| 888 | Wen, Milica Gašić, and Steve Young. 2017. <a href="#">PyDial:</a>            | <b>Dataset.</b> The client simulation subset contains                         | 943 |
| 889 | <a href="#">A Multi-domain Statistical Dialogue System Toolkit</a> .         | 2,176 instances across 28 client-side categories.                             | 944 |
| 890 | In <i>Proceedings of ACL 2017, System Demonstrations</i> ,                   | A train/test split of 1,604/572 utterances with                               | 945 |
| 891 | pages 73–78, Vancouver, Canada. Association for                              | conversation-level partitioning is used. The fol-                             | 946 |
| 892 | Computational Linguistics.   | lowing ablation studies (Sections A.2.2–A.2.3) are                            | 947 |
| 893 | Nicolas Wagner and Stefan Ultes. 2024. <a href="#">On the Con-</a>           | ducted on this subset to provide additional in-                               | 948 |
| 894 | <a href="#">trollability of Large Language Models for Dialogue</a>           | sights into design choices.   | 949 |
| 895 | <a href="#">Interaction</a> . In <i>Proceedings of the 25th Annual Meet-</i> |   |     |
| 896 | <i>ing of the Special Interest Group on Discourse and</i>                    |   |     |
| 897 | <i>Dialogue</i> , pages 216–221, Kyoto, Japan. Association                   |   |     |
| 898 | for Computational Linguistics.   |   |     |

| Hyperparameter                | Value                     |
|-------------------------------|---------------------------|
| <i>Training Configuration</i> |                           |
| Epochs                        | 10                        |
| Batch size                    | 64                        |
| Learning rate scheduler       | Cosine                    |
| Gradient clipping             | 1.0                       |
| Mixed precision (AMP)         | Enabled                   |
| <i>Early Stopping</i>         |                           |
| Patience                      | 3 epochs                  |
| Min delta                     | 0.001                     |
| Monitor metric                | Validation macro-F1       |
| <i>Model Architecture</i>     |                           |
| Context utterances            | {1, 4, 8, 12}             |
| History length                | 10 categories             |
| Attention heads               | 8                         |
| Transformer dropout           | 0.2                       |
| <i>Regularization Grid</i>    |                           |
| $\lambda_{tm}$                | {0.0, 0.2, 0.5, 1.0, 1.5} |
| Label smoothing $\epsilon$    | {0.0, 0.1, 0.2}           |
| <i>Cross-Validation</i>       |                           |
| Number of folds               | 5                         |
| Split strategy                | Conversation-level        |
| Random seed                   | 42                        |

Table 6: Hyperparameter settings. Context utterances and regularization weights were varied in grid search; other values were fixed across all experiments.

**Effect of Transition Regularization.** Table 7 shows the effect of varying  $\lambda_{tm}$  on the client simulation subset.

| $\lambda_{tm}$ | W-F1        | Top-3       | Trans.      | Cum70       |
|----------------|-------------|-------------|-------------|-------------|
| 0.0            | .225        | .451        | .599        | .801        |
| 0.2            | .255        | .495        | .743        | .894        |
| 0.5            | <b>.265</b> | <b>.505</b> | .799        | .925        |
| 1.0            | .251        | .497        | .838        | .955        |
| 1.5            | .253        | .494        | <b>.875</b> | <b>.961</b> |

Table 7: Effect of  $\lambda_{tm}$  on client simulation subset (28 categories). Trans.=transition validity. Cum70=cumulative accuracy at 70%.

TM regularization yields +18% weighted F1 improvement (0.225 $\rightarrow$ 0.265 at  $\lambda_{tm}$ =0.5), stronger than the +17% improvement on the full 60-category task. Transition validity improves by +46% (0.599 $\rightarrow$ 0.875 at  $\lambda_{tm}$ =1.5).

**Encoder Comparison.** Table 8 compares encoder performance on the client simulation subset.

Best weighted F1 (0.292) is achieved with GBERT-large at  $\lambda_{tm}$ =0.2. Notably, the absolute F1 scores on the 28-category subset are substantially higher than on the 60-category task (0.292 vs 0.164), reflecting the reduced complexity of the classification problem. However, the relative im-

| Encoder          | W-F1        | Top-3       | Params |
|------------------|-------------|-------------|--------|
| GBERT-large      | <b>.292</b> | .506        | 336M   |
| ModernGBERT-1B   | .287        | .507        | 1B     |
| GBERT-base       | .282        | <b>.518</b> | 110M   |
| EuroBERT-610M    | .281        | .492        | 610M   |
| EuroBERT-210M    | .274        | .513        | 210M   |
| ModernGBERT-134M | .271        | .538        | 134M   |
| GELECTRa-base    | .252        | .517        | 110M   |

Table 8: Encoder comparison on client simulation subset (28 categories, best  $\lambda_{tm}$  per encoder). Model size (110M–1B) shows no consistent correlation with performance.

provement from TM regularization remains consistent, confirming that dialogue-flow priors are valuable across task granularities.

### A.2.2 Transition-Matrix vs. Label Smoothing

Table 9 compares our transition-matrix regularization against label smoothing (Szegedy et al., 2015) with  $\epsilon \in \{0.0, 0.1, 0.2\}$  on the client simulation subset. While label smoothing provides modest gains for some architectures, transition-matrix regularization consistently outperforms across most encoder variants.

| Encoder          | Label Smoothing |                |                | TM Reg.     |          |
|------------------|-----------------|----------------|----------------|-------------|----------|
|                  | $\epsilon=0.0$  | $\epsilon=0.1$ | $\epsilon=0.2$ | Best        | $\Delta$ |
| GBERT-large      | .244            | .268           | .255           | <b>.292</b> | +0.024   |
| GBERT-base       | .253            | .263           | .260           | <b>.282</b> | +0.019   |
| ModernGBERT-134M | .256            | .251           | <b>.277</b>    | .271        | -.006    |
| GELECTRa-base    | .206            | .197           | .208           | <b>.253</b> | +0.044   |
| EuroBERT-210M    | .243            | .223           | .252           | <b>.275</b> | +0.023   |
| EuroBERT-610M    | .229            | .225           | .241           | <b>.281</b> | +0.041   |
| <i>Mean</i>      | .239            | .238           | .249           | <b>.276</b> | +0.024   |

Table 9: Label smoothing vs. TM regularization (weighted F1). TM outperforms on 5 of 6 encoders.

The mean improvement of TM regularization over the best label smoothing configuration ( $\epsilon=0.2$ ) is +0.024 F1 points. Notably, the gains are largest for architectures with lower baseline performance (GElectra-base: +4.4%, EuroBERT-610M: +4.1%), suggesting that domain-specific structural priors are particularly valuable when model capacity is limited. Table 4 in the main paper shows that this pattern holds even more strongly on the full 60-category task.

### A.2.3 Synthetic Data Augmentation

Given the substantial class imbalance in the client simulation subset (Gini coefficient 0.51), LLM-based synthetic data augmentation was explored.

A two-phase strategy was employed: (1) prompt-based augmentation using GPT-5-mini to balance minority classes up to 100 examples each, and (2) persona-based generation using expert-designed client profiles to introduce lexical and stylistic diversity. The augmented training set contains 8,487 instances (81% synthetic), reducing the Gini coefficient from 0.51 to 0.06.

Table 10 compares BERT-based models trained on real data only versus real+synthetic data on the client simulation subset.

| Encoder          | Real        | +Synth.     | $\Delta$ |
|------------------|-------------|-------------|----------|
| GBERT-large      | <b>.292</b> | .281        | -.011    |
| GBERT-base       | <b>.282</b> | .264        | -.018    |
| ModernGBERT-134M | .271        | <b>.283</b> | +.011    |
| GELECTRa-base    | .253        | <b>.257</b> | +.004    |
| EuroBERT-210M    | .275        | <b>.283</b> | +.008    |
| EuroBERT-610M    | <b>.281</b> | .276        | -.006    |
| <i>Mean</i>      | <b>.276</b> | .274        | -.002    |

Table 10: Real vs. real+synthetic training (weighted F1). Synthetic augmentation provides no consistent benefit for BERT-based models.

Results are mixed: synthetic data slightly improves smaller models (ModernGBERT-134M, EuroBERT-210M) but degrades larger models (GBERT-large, GBERT-base). The overall mean shows no significant difference ( $p > 0.17$ , independent t-test), suggesting that pretrained transformers are sufficiently data-efficient for this task.

In contrast, RNN baselines benefited substantially from synthetic augmentation: Simple RNN achieved F1=0.234 vs 0.097 on real data only. This suggests that pretrained transformers are sufficiently data-efficient for fine-grained dialogue act prediction, while RNNs require substantially more examples to converge. For this task, domain-specific structural priors (like transition matrices) appear more valuable than quantity-based augmentation when using pretrained encoders.

#### A.2.4 Cross-Dataset Validation on SWDA and Alternative KL Formulations

To test generalization beyond counselling and investigate the effect of skewed transition distributions, we evaluated TM regularization on the Switchboard Dialogue Act Corpus (SWDA; Shriberg et al., 1998; Stolcke et al., 2000), a benchmark of English telephone conversations filtered to 37 classes (excluding rare categories with fewer than 50 occurrences). Table 11 shows results using XLM-RoBERTa.

| $\lambda_{tm}$ | F1          | Top-3       | Cum70       | JS          |
|----------------|-------------|-------------|-------------|-------------|
| 0.0            | <b>.169</b> | .765        | .901        | .150        |
| 0.5            | .162        | .767        | .929        | .073        |
| 1.5            | .144        | <b>.771</b> | <b>.986</b> | <b>.023</b> |

Table 11: SWDA 37-class results. Dialogue-flow metrics (Cum70, JS) improve substantially but macro-F1 decreases with TM regularization.

Unlike counselling datasets (OnCoCo, HOPE), TM regularization on SWDA improves dialogue-flow alignment (JS: 0.150→0.023, Cum70: 0.901→0.986) but *decreases* macro-F1. This divergence reflects SWDA’s highly skewed transition distribution: analysis of the empirical transition matrix reveals that 69% of source classes have a single dominant successor—the “Statement-non-opinion” (sd) category, which accounts for 34% of all transitions. The regularizer thus pushes predictions toward this dominant class, improving flow alignment but suppressing minority class predictions.

To address this skewness, alternative formulations of the transition-matrix loss were explored: (1) reverse KL divergence, which is more permissive for confident predictions, and (2) entropy-weighted KL, which down-weights samples from high-entropy source categories where multiple successors are plausible. Neither variant improved over standard forward KL (Table 12), suggesting that the macro-F1 degradation on SWDA stems from the dataset’s inherent transition structure rather than the KL formulation.

| KL Variant            | Accuracy     | Macro-F1     |
|-----------------------|--------------|--------------|
| Forward KL (baseline) | <b>0.502</b> | <b>0.150</b> |
| Reverse KL            | 0.500        | 0.149        |
| Entropy-weighted      | 0.496        | 0.146        |
| Entropy + Reverse     | 0.495        | 0.146        |

Table 12: Alternative KL formulations on SWDA 37-class ( $\lambda_{tm}=0.5$ ). Forward KL performs best.

This finding suggests that TM regularization is most effective when transition patterns are more balanced, as in counselling domains where dialogue follows structured phases (problem exploration → intervention → resolution) rather than converging on a single dominant act type.

#### A.2.5 Context Window Size

The number of preceding utterances used as context is varied (1, 4, 8, 12). Table 13 shows macro-F1

across context lengths and TM weights. Performance peaks at 4 utterances with  $\lambda \geq 0.5$  (macro-F1=0.080). Longer context windows provide diminishing returns, suggesting that recent dialogue history is most informative for NDAP.

| Ctx | TM Weight ( $\lambda_{tm}$ ) |                 |                                   |                                   |                                   |
|-----|------------------------------|-----------------|-----------------------------------|-----------------------------------|-----------------------------------|
|     | 0.0                          | 0.2             | 0.5                               | 1.0                               | 1.5                               |
| 1   | .065 $\pm$ .012              | .069 $\pm$ .012 | .072 $\pm$ .010                   | .071 $\pm$ .011                   | .073 $\pm$ .009                   |
| 4   | .070 $\pm$ .019              | .073 $\pm$ .020 | <b>.080 <math>\pm</math> .013</b> | <b>.080 <math>\pm</math> .014</b> | <b>.080 <math>\pm</math> .011</b> |
| 8   | .066 $\pm$ .018              | .076 $\pm$ .015 | .074 $\pm$ .020                   | .079 $\pm$ .013                   | .079 $\pm$ .010                   |
| 12  | .065 $\pm$ .013              | .070 $\pm$ .019 | .075 $\pm$ .018                   | .077 $\pm$ .014                   | .077 $\pm$ .012                   |

Table 13: Macro-F1 by context window size and TM weight (mean  $\pm$  std, 5-fold CV). Best at 4 utterances with  $\lambda \geq 0.5$ .

### A.2.6 Hierarchical Label Embeddings

The OnCoCo taxonomy organizes 60 leaf categories into a five-level hierarchy. Experiments explored whether explicitly encoding this structure could improve predictions by embedding each hierarchy level ( $K_1$ – $K_5$ ) independently and integrating them with the conversation context via cross-attention. This approach is inspired by hierarchical text classification methods (Kowsari et al., 2017). However, experiments showed only marginal improvements over the base architecture (+0.5% weighted F1), insufficient to justify the additional complexity. The hypothesis is that the pretrained encoder already captures sufficient semantic structure, and that the transition-matrix regularizer—which implicitly encodes label relationships through observed co-occurrence patterns—provides a more effective inductive bias than explicit hierarchy embeddings.

### A.2.7 LLM Baseline Prompt Template

The following prompt template is used for GPT-5-mini zero-shot evaluation. Category descriptions and conversation history are inserted at the indicated placeholders.

You are an expert in dialogue act classification for German online counseling.

## Task

Predict the dialogue act category of the NEXT utterance in this counseling conversation.

## Categories (60 total)

[CATEGORY\_CODE]: [DESCRIPTION]

... (all 60 categories with descriptions)

## Conversation History (last 12 turns)

[SPEAKER] ([CATEGORY\_CODE] | [DESCRIPTION]): [TEXT]

... (conversation turns with speaker, category, text)

## Output Format (JSON)

Return your top 3 predictions:

```
{
  "predictions": [
    {"category": "CODE", "confidence": 0.6},
    {"category": "CODE", "confidence": 0.25},
    {"category": "CODE", "confidence": 0.15}
  ]
}
```

## B Significance Tests by Encoder

We report paired bootstrap significance tests (10,000 iterations) with Benjamini-Hochberg FDR correction for multiple comparisons. Table 14 summarizes results across regularization strengths, demonstrating consistent positive effects at all  $\lambda_{tm}$  values tested. Tables 15 and 16 show encoder-specific results.

| $\lambda_{tm}$ | Tests | Pos. | %Pos. | Sig. | %Sig. | Mean $\Delta$ | Med. $\Delta$ |
|----------------|-------|------|-------|------|-------|---------------|---------------|
| 0.2            | 42    | 36   | 86%   | 9    | 21%   | +0.0061       | +0.0050       |
| 0.5            | 44    | 37   | 84%   | 25   | 57%   | +0.0069       | +0.0063       |
| 1.0            | 45    | 40   | 89%   | 22   | 49%   | +0.0089       | +0.0091       |
| 1.5            | 42    | 33   | 79%   | 19   | 45%   | +0.0075       | +0.0070       |

Table 14: Summary of TM regularization significance tests by  $\lambda_{tm}$  value. Tests: paired comparisons (encoder  $\times$  architecture  $\times$  context). Pos.: positive effect over  $\lambda_{tm}=0$ . Sig.: significant after FDR correction ( $\alpha=0.05$ ).

| Encoder          | $\Delta$ Macro-F1 | Sig. (FDR<0.05)    |
|------------------|-------------------|--------------------|
| ModernGBERT-134M | +0.93%            | 5/8 (62%)          |
| EuroBERT-610M    | +0.80%            | 4/5 (80%)          |
| ModernGBERT-1B   | +0.72%            | 3/7 (43%)          |
| GBERT-base       | +0.70%            | 5/6 (83%)          |
| GELECTRa-base    | +0.68%            | 4/8 (50%)          |
| GBERT-large      | +0.47%            | 3/7 (43%)          |
| EuroBERT-210M    | +0.38%            | 1/3 (33%)          |
| <b>Overall</b>   | <b>+0.69%</b>     | <b>25/44 (57%)</b> |

Table 15: TM effect ( $\lambda_{tm}=0$  vs 0.5) by encoder. Mid-sized encoders show highest significance rates.

| Encoder          | $\Delta$ Macro-F1 | Sig. (FDR<0.05)    |
|------------------|-------------------|--------------------|
| GELECTRa-base    | +1.93%            | 4/4 (100%)         |
| ModernGBERT-134M | +0.93%            | 2/4 (50%)          |
| EuroBERT-210M    | +0.75%            | 2/3 (67%)          |
| EuroBERT-610M    | +0.64%            | 2/4 (50%)          |
| ModernGBERT-1B   | +0.50%            | 1/4 (25%)          |
| GBERT-large      | −0.08%            | 1/3 (33%)          |
| GBERT-base       | −0.16%            | 1/4 (25%)          |
| <b>Overall</b>   | <b>+0.67%</b>     | <b>13/26 (50%)</b> |

Table 16: Architecture effect (BERT $\rightarrow$ History,  $\lambda_{tm}=0.5$ ) by encoder. Benefits weaker encoders only.