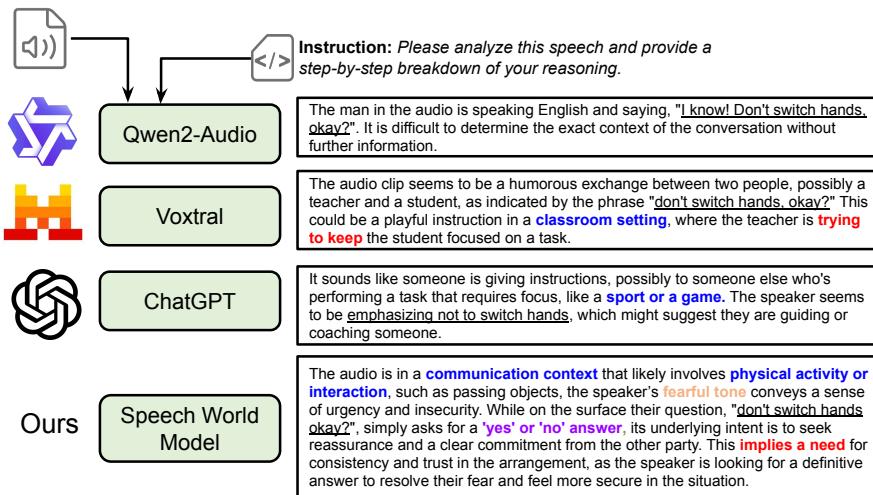


000 SPEECH WORLD MODEL: CAUSAL STATE–ACTION 001 PLANNING WITH EXPLICIT REASONING FOR SPEECH 002

003 **Anonymous authors**
004

005 Paper under double-blind review
006



007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Figure 1: Speech World Model. Demo audio link: <http://bit.ly/4pBJuWP>.

ABSTRACT

Current speech-language models (SLMs) typically use a cascade of speech encoder and large language model, treating speech understanding as a single black box. They analyze the content of speech well but reason weakly about other aspects, especially under sparse supervision. Thus, we argue for explicit reasoning over speech states and actions with modular and transparent decisions. Inspired by cognitive science we adopt a modular perspective and a world model view in which the system learns forward dynamics over latent states. We factorize speech understanding into four modules that communicate through a causal graph, establishing a cognitive state search space. Guided by posterior traces from this space, an instruction-tuned language model produces a concise causal analysis and a user-facing response, enabling counterfactual interventions and interpretability under partial supervision. We present the first graph based modular speech model for explicit reasoning and we will open source the model and data to promote the development of advanced speech understanding.

1 INTRODUCTION

Speech language models (SLMs) (Cui et al., 2024) provide a unified framework for general-purpose speech understanding, enabling multitask perception, reasoning, and multi-turn interaction. At the same time, an ongoing debate persists regarding whether such models truly exhibit genuine reasoning capabilities or merely perform sophisticated pattern matching and statistical abstraction over tokens (Shojaee et al., 2025). One direct observation is that current SLMs appear to combine outputs from isolated tasks—such as automatic speech recognition, emotion recognition, communicative function detection, intent classification, and pragmatic inference—and then present this aggregation *as if it reflects reasoning*. However, this overlooks the *intrinsic dependencies among these speech components* and provides only partial supervision, which we argue is fundamental to enabling genuine

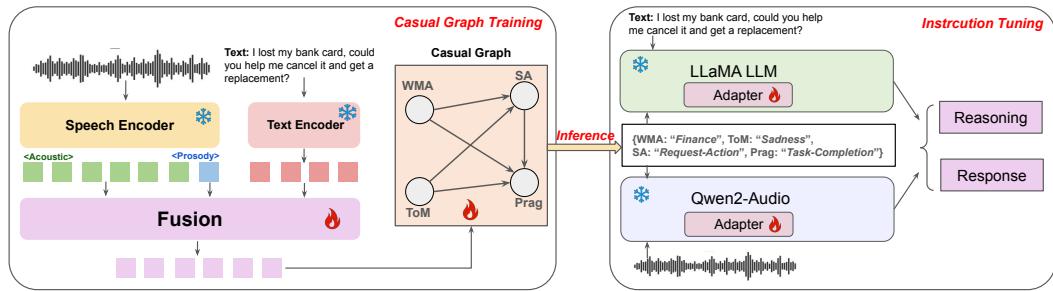


Figure 2: The pipeline of our *Speech World Model*, a two-stage training framework we term “Causal Graph-Guided Explicit Reasoning”. In the first stage, a Causal Graph is trained to convert speech into latent states as structured reasoning information. In the second stage, this information is used as an explicit guide for the parameter-efficient instruction tuning of large (speech) language models to enhance their speech understanding and reasoning capabilities.

reasoning in speech understanding. From a human-centered cognitive perspective, speech perception is thought to be hierarchically organized into partially specialized modules—such as acoustic, linguistic, and paralinguistic processing—that interact through causal and reciprocal relations (Hickok & Poeppel, 2007). However, this does not entail designing strictly brain-like modular systems (e.g. verified AI (Seshia et al., 2022)), as human perception itself operates through overlapping and partially opaque processes that remain incompletely understood (Angelaki et al., 2025).

An alternative line of research seeks an intermediate approach that models latent internal states, thereby enabling more interpretable predictions of system dynamics. A representative example is *World Models* (Ha & Schmidhuber, 2018; LeCun, 2022), which posit that the subsequent state is conditionally predictable from any given current state. In the domain of speech, human perceptual systems also anticipate interdependent dynamics across states such as emotion, intent, and linguistic content, such that changes in one dimension systematically constrain others (Hickok & Poeppel, 2007). These observations point to an underlying structure of causal dependencies. It is worth noting that Chain-of-Thought (CoT) prompting (Wei et al., 2022b) represents a parallel approach, which expands the search space to improve goal-directed reasoning at the expense of increased computational cost. However, the search process is not grounded in human speech perception principles.

In this work, we propose *Speech World Models* (SWM) to advance speech understanding and reasoning by modeling the dynamics of internal speech states. SWM comprises four explicit modules—*World Model Action* (*WMA*) (Mesaros et al., 2018), *Theory of Mind* (*ToM*) (Sclar et al., 2025), *Speech Act* (*SA*) (Bothe et al., 2020), and *Pragmatic Intent* (*Prag*) (Stolcke et al., 2000)—as an approximated subspace of the cognitive perception system. A causal-graph-based perception module is introduced to bridge these components and capture the flow of causal information. In the first stage, this graph is trained in both supervised and semi-supervised settings to establish a cognitive state search space. Notably, leveraging causal dependencies among modules yields faster convergence than training the modules independently. Moreover, because the graph approximates cognitive state structure, it can reliably infer unlabeled modules from partially labeled data, suggesting that causal graphs function as natural generative structures for latent variables. Once the search space is established, in the second stage we incorporate the outputs of the causal graph—samples from the search space, state nodes, and information flow—into prompts for instruction tuning. This provides an explicit reasoning advantage: by enforcing structured reasoning chains, SLMs are guided to reduce hallucinations and achieve greater overall reasoning capacity.

Experiments demonstrate the superiority of SWM over traditional SLMs across a range of reasoning metrics. First, when comparing against a randomly initialized graph without cognitive grounding, we find that our causal graph converges faster during training and achieves substantially higher reasoning scores (e.g., ACE, ICS). Second, instruction-tuning results show that SWM outperforms open-source SLMs such as Qwen1 (Chu et al., 2023), Qwen2-Audio (Chu et al., 2024), and Voxtral (Liu et al., 2025) in reasoning ability, with particularly strong performance in emotion recognition—even surpassing commercial models—highlighting the benefits of explicit reasoning. Third, while the overall speech understanding performance of SWM (as measured by our proposed metrics) remains

108 slightly below that of Gemini 2.5 Pro (Gemini Team, 2025), our training cost is dramatically lower
 109 (only 20 GPU hours), validating the efficiency of the causal-graph-based explicit reasoning paradigm.
 110

111 2 RELATED STUDY

113 2.1 COGNITIVE FOUNDATIONS OF SPEECH PROCESSING

115 Cognitive neuroscience shows that speech perception and comprehension emerge from modular
 116 computations over latent states, rather than a single monolithic process. Core frameworks include
 117 the dual-stream model linking ventral (semantic) and dorsal (sensorimotor) pathways (Hickok &
 118 Poeppel, 2007), staged accounts of syntactic and semantic processing (Friederici, 2011), and analyses
 119 of the superior temporal gyrus that reveal intermediate phonological and prosodic representations
 120 (Leonard & Chang, 2022). Complementary perspectives highlight cortical oscillations aligned with
 121 hierarchical linguistic units (Giraud & Poeppel, 2012) and predictive coding views of the brain as a
 122 generative model updating latent states via prediction errors (Friston et al., 2021). Together, these
 123 perspectives converge on four principles: (i) modular and hierarchical decomposition of speech, (ii)
 124 forward predictive dynamics, (iii) integration of perception and action, and (iv) generative modeling
 125 of latent auditory states. Our Speech World Model operationalizes these principles computationally.

126 2.2 SPEECH LANGUAGE MODELS AND REASONING

127 Early work on speech language models (Gong et al., 2024; 2023; Tang et al., 2024; Kong et al.,
 128 2024) pioneered the use of instruction tuning pipelines for speech understanding and preliminary
 129 reasoning tasks. Subsequent studies adopted Chain-of-Thought (CoT) prompting (Wei et al., 2022b)
 130 to construct reasoning-chain search spaces, which advanced reasoning capabilities in speech-related
 131 tasks (Ma et al., 2025; Xie et al., 2025a;b; Kong et al., 2025); see also (Cui et al., 2024) for a broader
 132 review. Nevertheless, these CoT-based approaches remain largely heuristic and are not grounded in
 133 principles of human speech cognition (Hickok & Poeppel, 2007).

134 2.3 WORLD MODELS

135 The concept of world models has long-standing roots in cognitive science and neuroscience, where
 136 internal models are thought to support prediction, planning, and reasoning in human perception and
 137 action (Craik, 1967; Friston, 2005; Hickok & Poeppel, 2007). Ha & Schmidhuber (2018) popularized
 138 the term by introducing a generative framework based on latent dynamics, demonstrating how agents
 139 could “dream” in an internal environment. This line of work inspired subsequent advances such
 140 as PlaNet (Hafner et al., 2019b), Dreamer (Hafner et al., 2019a), and MuZero (Schrittwieser et al.,
 141 2020), which emphasized model-based reinforcement learning and planning from learned dynamics.
 142 More recently, LeCun (2022) argued that world models should be central to autonomous machine
 143 intelligence, framing them as self-supervised systems that learn hierarchical predictive representations.
 144 Parallel efforts have explored domain-specific instantiations for vision (Garrido et al., 2024), motor
 145 control (Assran et al., 2025), and language tasks (Lin et al., 2024; Feng et al., 2025).

146 3 SPEECH WORLD MODEL SYSTEMS

147 In this section, we present the methodology behind our proposed *Speech World Model (SWM)*. At
 148 its core, the framework relies on a causal graph that factorizes speech understanding into modular
 149 components, enabling explicit reasoning over speech states and actions. These states then condition
 150 (spoken) language models, which are fine-tuned to generate both a transparent reasoning trace
 151 and a potential response to given utterance. This design allows us to validate the hypothesis that
 152 explicit state representation improves Speech Language Model’s reasoning capacity and reduces
 153 hallucinations (Atwany et al., 2025).

154 3.1 WORLD MODEL CAUSAL GRAPH

155 We proposed a probabilistic causal model, represented by a Directed Acyclic Graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$,
 156 where \mathcal{V} is a set of nodes representing different modules and \mathcal{E} is a set of directed edges representing
 157 causal relationships between modules. This structure allows us to explicitly model the dependencies
 158 between different aspects of speech and perform causal reasoning.

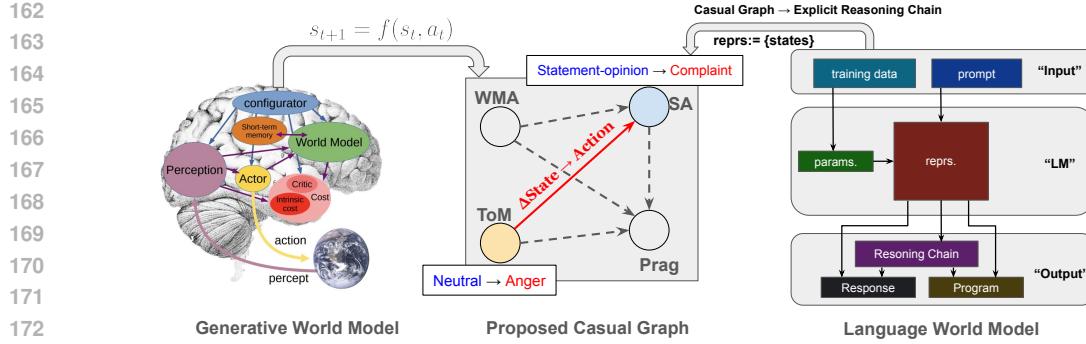


Figure 3: A unified perspective on world models. Both the Generative (left(Garrido et al., 2024)) and Language (right) world models are forward dynamic models. Our Causal Graph (center) provides a structured, explicit formulation of these dynamics. See Appendix. A.1 for details.

3.1.1 MODULES AND STATES

We define four key modules, each capturing a distinct aspect of speech understanding, forming the *cognitive perceptual subspace*. The complete label spaces are provided in the Appendix. A.4.2.

- **World Model Activation (WMA)** (Mesaros et al., 2018) captures context and situation of the speech environment.
- **Theory of Mind (ToM)** (Sclar et al., 2025) models the speaker’s mental and emotional states.
- **Speech Act (SA)** (Bothe et al., 2020) identifies the communicative function of an utterance.
- **Pragmatic Intent (Prag)** Stolcke et al. (2000) represents the speaker’s underlying intentions.

Graph Perspective. Each module corresponds to a node, instantiated as a neural network performing its own computation. The outputs of these modules are categorical latent states, and collectively they form the structured representation of speech, as presented in Fig. 3.

World Model Perspective. Inspired by the essence of image world models (Garrido et al., 2024), we view the entire system as a forward dynamics model: the current latent state, together with an action, leads to the next latent state. Within this view, the notion of action extends beyond external commands. Instead, an action can be understood as the causal influence exerted by one state on another. For example, when the Theory of Mind module shifts from *neutral* to *anger*, the downstream Speech Act may change from *Statement-opinion* to *Complaint/Escalate*, as shown in Fig. 3.

Thus, actions in our framework are best understood as the state transitions induced by causal relationships within the graph. Formally, for a node v with parents $\text{Pa}(v)$:

$$S_v = f_v(\{S_u : u \in \text{Pa}(v)\}, A_{u \rightarrow v}), \quad (1)$$

where S_v is the latent state at node v , and $A_{u \rightarrow v}$ denotes the action, interpreted as the causal influence (action) of parent node u on v .

3.1.2 GRAPH CALCULATION

Given an input speech signal X , the causal graph infers a joint posterior distribution over latent states $Z = \{z_{WMA}, z_{ToM}, z_{SA}, z_{Prag}\}$. By the DAG factorization:

$$\begin{aligned} p(Z|X) &= p(z_{WMA}|X) \cdot p(z_{ToM}|X) \cdot \\ &\quad p(z_{SA}|z_{WMA}, z_{ToM}, X) \cdot p(z_{Prag}|z_{SA}, z_{ToM}, z_{WMA}, X) \end{aligned} \quad (2)$$

Each conditional $p(z_v|\text{Pa}(v), X)$ is parameterized by a neural module $f_v(\cdot; \theta_v)$, modeled as a multi-class classifier. Our proposed causal graph is indicated in Fig. 4 (B). The detailed calculation and model architecture are detailed in Appendix. A.2.1.

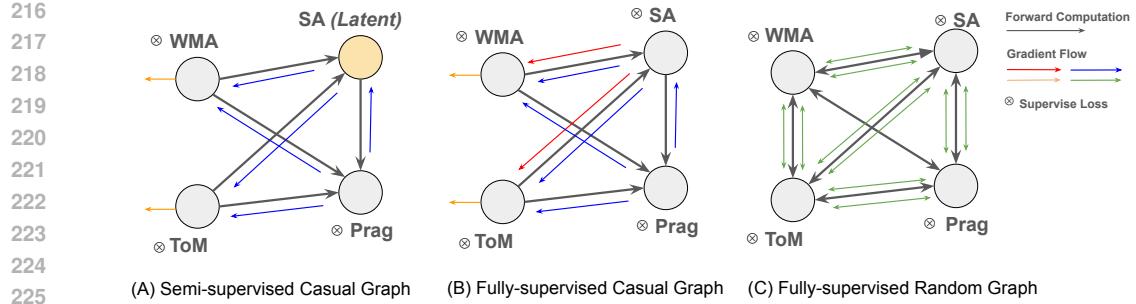


Figure 4: Comparison of gradient flow under different training scenarios. (A) Semi-supervised causal graph: when a module lacks direct supervision, it act as latent variable generator, and gradients propagate backward through causal edges from downstream modules. (B) Fully-supervised causal graph: all modules have labels, so losses are applied locally, while causal dependencies still guide gradient flow. (C) Fully-supervised random graph: supervision is present at every node but edges are unstructured, leading to redundant and less efficient gradient propagation.

3.1.3 MULTI-TASK TRAINING WITH CASUAL GRAPH

Causal graph is trained in a multi-task manner to form reasoning-chain search space. The supervised loss is shown in Eq.3, where $m_{i,v}$ indicates whether the label for node v is available in sample i .

$$\mathcal{L}_{\text{sup}} = \sum_{i=1}^N \sum_{v \in \mathcal{V}} m_{i,v} \text{CE}(y_{i,v}, S_{i,v}), \quad (3)$$

Teacher forcing is applied edge-wise: each child node receives either the ground-truth parent state or the predicted distribution. For edge $u \rightarrow v$:

$$\tilde{S}_{i,u} = \tau_{i,u \rightarrow v} \text{onehot}(y_{i,u}) + (1 - \tau_{i,u \rightarrow v}) \text{stopgrad}(S_{i,u}), \quad \tau_{i,u \rightarrow v} \sim \text{Bernoulli}(p_{u \rightarrow v}). \quad (4)$$

The mixed signal $\tilde{S}_{i,u}$ is used as the parent input when computing the state of child v .

3.1.4 SEMI-SUPERVISED LEARNING

Motivation. Annotations across modules are costly and often incomplete (e.g., WMA labels are missing more frequently). If parameter updates were restricted only to nodes with available labels, a large fraction of training data would be wasted. Semi-supervised training allows unlabeled parent nodes to be optimized through the supervision signals of their labeled children.

Gradient Flow. Consider a parent node j without labels in sample i ($m_{i,j} = 0$), where one of its child node k has a label ($m_{i,k} = 1$). To ensure a differentiable path to the unlabeled parent node j , teacher forcing is disabled ($p_{j \rightarrow k} = 0$), forcing the child k to rely on the parent's continuous prediction $S_{i,j}$. Thus, the loss $\mathcal{L}_{i,k} = \text{CE}(y_{i,k}, S_{i,k})$ propagates to θ_j through the chain rule:

$$\frac{\partial \mathcal{L}_i}{\partial \theta_j} = \sum_{\substack{k: m_{i,k}=1 \\ j \in \text{Pa}(k)}} \frac{\partial \mathcal{L}_{i,k}}{\partial \eta_{i,k}} \frac{\partial \eta_{i,k}}{\partial S_{i,j}} \frac{\partial S_{i,j}}{\partial \eta_{i,j}} \frac{\partial \eta_{i,j}}{\partial \theta_j} \neq 0. \quad (5)$$

where the gradient also flows into the parent distributions $p(s_u | \xi_u)$. Thus unlabeled parents act as *latent variable generators* that receive updates from their supervised children as shown in Fig. 4 (A).

Multi-task supervision ensures “learn whenever a label is available” (Eq. 3), while semi-supervised training ensures “learn even when labels are missing” via chain gradients (Eq. 5). Together they allow effective learning from incompletely annotated data under the causal graph structure.

3.2 RANDOM GRAPH

To validate the benefits of our proposed causal structure, grounded by cognitive speech perception, we introduce the *Random Graph* as a strong, unstructured and data-driven baseline. This model replaces

270 the predefined DAG with a fully connected graph where every module can, in principle, influence
 271 every other module. Consequently, the prediction for each module is conditioned on both the speech
 272 input X and the initial, teacher-forced estimates from **all** other modules in the graph, allowing the
 273 model to learn arbitrary dependencies directly from the data.

$$275 \quad p(Z|X) = \prod_{v \in \mathcal{V}} p(z_v|Z \setminus \{z_v\}, X) \quad (6)$$

277 By comparing this model with our structured causal graph, we can quantify the contribution of
 278 imposing theoretically-grounded causal priors on speech understanding. The detailed, step-by-step
 279 computation for this iterative process is provided in Appendix. A.2.2.

281 3.3 SPEECH UNDERSTANDING AND REASONING

283 To bridge the gap between the structured latent states inferred by the causal graph and human-
 284 interpretable analysis, we employ an instruction tuning (Wei et al., 2022a) phase. This final stage
 285 is designed to validate our core hypothesis: *Conditioning on these explicit graph states tends to*
 286 *guide the reasoning-chain search toward human-aligned spaces..* To this end, the model is trained
 287 to produce a transparent reasoning process and a contextually appropriate response based on the
 288 graph's outputs. We explore two distinct settings for this purpose: a language-only approach and a
 289 multi-modal approach.

290 3.3.1 LANGUAGE-ONLY SETTING

292 In this setting, we fine-tune a standard LLM to reason over the symbolic outputs of the pre-trained
 293 causal graph \mathcal{G} . The input is prompted with a general instruction (Instr) and a textual serialization of
 294 the graph's output, $I(\mathcal{G}(x))$, which encodes the inferred states and their causal relationships. The
 295 prompts we used for training are detailed in Appendix. A.5.2.

296 The model is trained to produce a target sequence y that contains both a step-by-step *reasoning*
 297 *process* and a potential *dialogue response*. The training objective is to minimize the standard cross-
 298 entropy loss over the target sequence, where θ represents the parameters of the LLM, and the target y
 299 is the ground-truth reasoning and response pair.

$$301 \quad \mathcal{L}_{\text{IT}}(\theta) = - \sum_{(x,y) \in \mathcal{D}} \log p_\theta(y \mid \text{Instr}, I(\mathcal{G}(x))) \quad (7)$$

304 3.3.2 MULTI-MODAL SETTING

306 This approach leverages a speech language model to directly ground the symbolic states in the
 307 acoustic properties of the raw speech signal. The input to the SLM consists of the instruction (Instr),
 308 the raw speech input x , and the set of inferred states $\{S_v\}_{v \in \mathcal{V}}$. The model's task is to generate the
 309 same target sequence y containing a reasoning process and a response.

$$310 \quad \mathcal{L}_{\text{IT}}(\theta) = - \sum_{(x,y) \in \mathcal{D}} \log p_\theta(y \mid \text{Instr}, x, S_{\text{WMA}}, S_{\text{ToM}}, S_{\text{SA}}, S_{\text{Prag}}) \quad (8)$$

313 This objective forces the SLM to build a joint understanding of the audio signal and the structured
 314 states, effectively learning to associate “what was said” and “how it was said” with the “why” provided
 315 by our causal graph’s analysis.

316 4 THEORETICAL ANALYSIS OF THE SPEECH WORLD MODEL

318 The use of a causal graph for speech understanding leads to more efficient training by providing a
 319 strong structural prior. This prior improves learning in two key ways: by optimizing gradient flow
 320 within the graph and by constraining the search space for downstream instruction tuning. See more
 321 detailed derivations in Appendix. A.3.

323 **Structured Priors and Optimized Gradient Flow.** The causal graph simplifies the learning
 324 problem by modeling the joint distribution over latent variables as a factorized probability, where

324 each module is conditioned only on its causal parents (Pearl, 2009):
 325

$$P(Y_1, \dots, Y_N | X) = \prod_{i=1}^N P(Y_i | X, Y_{\text{Pa}(i)}) \quad (9)$$

326 This factorization directly yields a more efficient and stable gradient flow. The gradient for a module v
 327 is computed based only on its parents, precluding the need to discover relationships from scratch and
 328 focusing updates along causally relevant pathways. This structure minimizes *learning redundancy*
 329 defined as wasteful gradient updates across unrelated components, which we can quantify as:
 330

$$\text{Redundancy} \propto \sum_{i=1}^N \Pi_i \|\nabla_{\theta_i} \mathcal{L}(\theta_i)\| \quad (10)$$

331 where Π_i is an importance indicator for module i . By pruning unnecessary dependencies, the
 332 model’s transparent design directly accelerates convergence, a known benefit of causality-inspired
 333 explainability (Schölkopf et al., 2021).
 334

335 **Constrained Search Space for Instruction Tuning.** This causal structure provides a critical
 336 advantage for downstream instruction tuning. The graph constrains the combinatorially large search
 337 space of latent cognitive states, $\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_N$, to a valid, low-entropy subspace $\mathcal{S}_{\mathcal{G}} \subset \mathcal{Y}$. This
 338 reduction in complexity is quantifiable by information entropy:
 339

$$H(Y | X, \mathcal{G}) = \sum_{i=1}^N H(Y_i | Y_{\text{Pa}(i)}, X, \mathcal{G}) \ll \sum_{i=1}^N \log |\mathcal{Y}_i| \quad (11)$$

340 This simplification transforms the learning task for the language model. Instead of learning a complex
 341 joint policy $\pi_{\phi}(Y, R | X)$ to generate both a reasoning chain Y and a response R , as in standard Chain-
 342 of-Thought (Wei et al., 2022b; Yao et al., 2023), our method allows learning a simpler conditional
 343 policy $\pi_{\phi}(R | X, Y_{\mathcal{G}})$. Providing the causally valid state $Y_{\mathcal{G}}$ as input acts as a *control variate*, reducing
 344 gradient variance and significantly improving the sample efficiency of instruction tuning.
 345

351 5 EXPERIMENTS

352 5.1 DATA

354 5.1.1 REAL-WORLD DATASET

- 356 • **MELD** (Poria et al., 2019) is a conversational dataset for emotion recognition. It consists of over
 357 13,000 utterances from the TV series Friends, annotated with emotion and sentiment labels.
- 358 • **IEMOCAP** (Busso et al., 2008) is an audiovisual and multi-speaker database of dyadic interactions
 359 between actors, containing approximately 12 hours of data annotated with categorical and
 360 dimensional emotion labels.
- 361 • **SLURP** (Bastianelli et al., 2020) is a challenging spoken language understanding (SLU) dataset
 362 featuring single-turn user interactions with a voice assistant. Utterances are annotated with a
 363 hierarchical schema, providing both user intent and corresponding action.
- 364 • **VoxCeleb** (Nagrani et al., 2017) is a large-scale text-independent speaker identification dataset. It
 365 comprises hundreds of thousands of utterances from celebrities, extracted from YouTube videos.

366 We use the above datasets as the primary speech sources for our experiments. Detailed statistics of
 367 these datasets are provided in Appendix A.4.1.

369 5.1.2 LABEL-GENERATION

370 A central challenge is that speech data usually comes with only *partial supervision*: some modules
 371 may be annotated while others are not. To construct complete instances for fully-supervised graph
 372 training and instruction tuning, we design a two-stage label generation procedure.
 373

374 **Stage 1: Label completion.** Given transcription and any known module labels, the goal is to predict
 375 missing labels for the remaining modules. We prompt a large language model Vicuna (Chiang et al.,
 376 2023) with the complete label space and enforce strict decoding into valid categories. This stage can
 377 be interpreted as a constrained imputation problem, where fuzzy mapping is applied to handle minor
 378 mismatches between predicted tokens and the canonical label set.

378 **Stage 2: Reasoning and response generation.** With the completed labels and transcription as
 379 input, we query the LLM again to generate two additional outputs: (i) a comprehensive reasoning
 380 analysis how the predicted states interact through the causal graph and jointly shape the interpretation
 381 of the utterance, and (ii) a potential response aligned with the inferred states.
 382

383 Full details of prompts, model configuration, and other implementation specifics are provided in the
 384 Appendix. A.4.2.

385 5.2 METRICS

387 **Node quality.** We evaluate the performance of individual modules (nodes) on their respective
 388 sub-tasks using standard classification *Accuracy* and *weighted F1-Score*.

389 **Edge causal effect.** To validate each learned edge $X \rightarrow Y$, we measure its causal strength using
 390 the *Average Causal Effect (ACE)* (Pearl, 2010) and its directional consistency with the *Intervention*
 391 *Consistency Score (ICS)*. *ACE* indicates the magnitude of post-intervention effects, while *ICS* indicates
 392 if effects align with a predefined hypothesis. See Appendix. A.8.1 for computational details.

393 **Instruction tuning.** To provide a scalable and consistent evaluation of our instruction-tuned models,
 394 we employ the Model-as-Judge (M.J.) methodology (Zheng et al., 2023), where GPT-4o (OpenAI,
 395 2024) acts as an impartial judge to score both reasoning and response quality. The overall M.J. score
 396 is a weighted sum of reasoning and response scores. *EM* and *EA* refer to emotion mention rate and
 397 emotion classification accuracy, respectively, while *R-Len* represents the length of the reasoning
 398 process. The detailed prompting strategy and scoring rubric are indicated in Appendix. A.9.

399 5.3 EXPERIMENTS SETUP

401 Our method involves a two-stage training process: causal graph training followed by instruction
 402 tuning. We train the causal graph under three settings: fully-supervised (for causal and random
 403 graphs) and semi-supervised (for the causal graph). Each module is an MLP-based classifier, with the
 404 WMA and ToM modules enhanced by a preceding temporal attention layer. We also conduct ablation
 405 studies on the causal graph, with detailed results available in Appendix. A.6. For instruction tuning,
 406 we use LoRA (Hu et al., 2022) for parameter-efficient instruction tuning of Llama3.1-8B Meta (2024)
 407 (language-only) and Qwen2-Audio (Chu et al., 2024) (multi-modal). As a baseline, we also fine-tune
 408 Qwen2-Audio using only CoT-style prompts, excluding any outputs from our causal graph. Detailed
 409 model architecture training configurations can be found in Appendix. A.5 and A.7.

410 5.4 GRAPH EVALUATION

412 The results of our graph evaluation are presented in Table. 1 and Table. 2, with a detailed breakdown
 413 of edge causal effects shown in Fig. 5. These results first demonstrate the **efficiency** of our Causal
 414 Graph, which **achieves node accuracy** comparable to a Random Graph baseline at a significantly
 415 **lower training cost**. More importantly, it provides **stable and interpretable causality** by quantifying
 416 effects with *ACE* and *ICS* metrics. In contrast, the Random Graph’s internal logic is proven unstable
 417 and dependent on training hyperparameters. Furthermore, our Causal Graph proves robust under
 418 incomplete supervision. In semi-supervised settings, it **maintains strong performance**, with any
 419 reduction in causal strength logically localized to the edges connected to the unsupervised module,
 420 highlighting the reliability of its structure.

421 Table 1: Node and Edge performance of the Causal Graph under various supervision settings.

Method	Setting	Node Quality (Accuracy %, ↑)				Edge Casual Effect	
		WMA	ToM	SA	Prag	Ave. ACE (%), ↑	Ave. ICS (%), ↑
Fully-supervised	–	69.4	73.5	65.3	81.4	23.57	43.29
Semi-supervised	Latent module: WMA	34.8	75.0	70.7	83.2	21.71	26.9
	Latent module: ToM	69.1	43.3	69.6	83.5	21.98	28.9
	Latent module: SA	69.3	77.0	34.4	82.5	21.65	29.3
Random Graph	–	69.7	74.0	67.5	83.6	–	–

430 5.5 SPEECH UNDERSTANDING AND REASONING EVALUATION

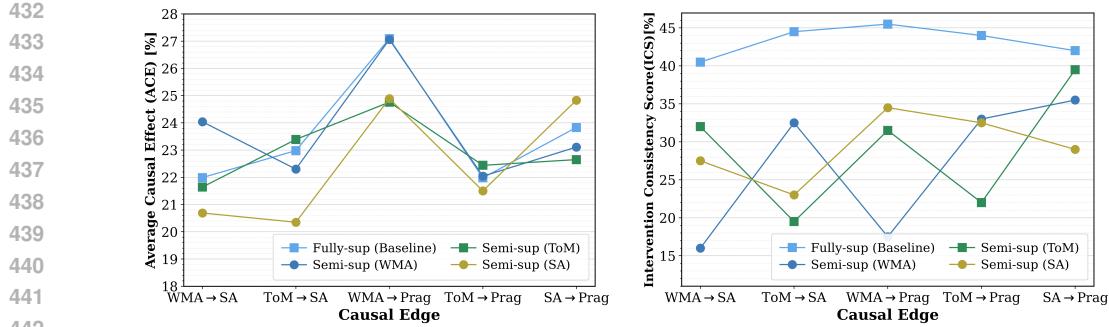


Figure 5: ACE and ICS of each causal edge under both fully-supervised and semi-supervised training.

Our speech understanding and reasoning results (Table. 3) reveal the distinct contributions of our reasoning framework. To isolate these effects, we created a tuned baseline by fine-tuning Qwen2-Audio with our CoT-style instruction data. This **baseline alone surpasses other open-source models**, confirming the high quality of our instruction tuning data, and our full SWM models, which leverage the causal graph for explicit reasoning, **significantly outperform this tuned baseline**. This demonstrates that while good data is important, the causal graph’s explicit guidance is the critical factor driving superior reasoning ability. This benefit is particularly stark in **emotion recognition**, where our model’s accuracy exceeds even top proprietary models. We posit that this high fidelity stems from the graph grounding the model’s analysis, thus reducing hallucination and enforcing a more structured reasoning process. Finally, while a model like Gemini 2.5 Pro (Gemini Team, 2025) holds an edge in the overall score, **our approach achieves this competitive performance at a dramatically lower training cost**. This validates our causal graph-guided paradigm as a highly efficient and effective path to complex speech understanding, prioritizing resourcefulness without a significant compromise in quality.

Table 2: Node quality and information flow analysis (see Appendix.A.8.2for details) for random graph baseline.

$p(\text{teacher_force})$	Node Quality (Accuracy %, ↑)			Information-flow	
	WMA	ToM	SA	Strongest	Weakest
0.0	70.1	75.2	66.1	82.8	SA → WMA
0.3	69.7	74.0	67.5	83.6	ToM → SA
0.5	69.5	73.1	65.9	82.0	WMA → SA
0.8	69.4	74.4	70.3	83.2	ToM → Prag
1.0	70.1	73.9	66.3	82.3	Prag → ToM
				SA → ToM	ToM → Prag

Table 3: Performance comparison against open-source and proprietary baselines.

Method	Prompt Style (Inference)	Overall M.J. Score ↑ $0.6 \times R_s + 0.4 \times R_p$	Reasoning Score ↑ R_s	Response Score ↑ R_p	Reasoning Breakdown (%) ↑ EM EA	R-Len (words)
<i>Our Models</i>						
SWM (Llama3.1-8b)	CoT	7.81	7.84	7.76	97.80	66.26
SWM (Qwen2-Audio)	CoT	7.59	7.26	8.08	91.80	71.02
<i>Tuned Baseline</i>						
Qwen2-Audio-CoT	CoT	5.18	4.76	5.82	92.11	34.72
<i>Baselines</i>						
Qwen-Audio (Chu et al., 2023)	Direct	2.70	2.20	3.46	14.20	8.00
Qwen2-Audio (Chu et al., 2024)	Direct	2.63	2.08	3.47	5.14	15.38
Qwen2-Audio (Chu et al., 2024)	CoT	2.39	1.96	3.04	6.11	17.50
Voxtral (Liu et al., 2025)	Direct	2.89	2.46	3.54	10.28	5.88
Voxtral (Liu et al., 2025)	CoT	2.92	2.52	3.52	10.89	5.56
<i>Proprietary Models</i>						
GPT-4o (OpenAI, 2024)	CoT	7.41	6.98	8.06	68.20	45.16
Gemini 2.5 Pro (Gemini Team, 2025)	CoT	8.12	8.02	8.28	82.47	51.29

6 LIMITATIONS AND CONCLUSIONS

In this work, we proposed Speech World Model (SWM), which surpasses the current open-sourced speech language models in understanding and reasoning ability by modeling the dynamics of internal speech states. However, there are still some limitations. First, the current model includes only four modules, while more diverse states could enhance speech understanding and better capture complex speech dynamics. Second, the causal graph in SWM is predefined, limiting the model’s ability to adapt to unseen dependencies. Future work could focus on learning adaptive causal structures to increase flexibility. Third, while instruction tuning effectively integrates reasoning traces and responses, it relies heavily on the label generation pipeline, where inaccuracies can propagate to both the reasoning and response stages. Therefore, achieving efficient data annotation and generation for speech remains a crucial and collaborative challenge within the speech community.

486 **7 ETHICS STATEMENT**
 487

488 Our work adheres to the ICLR Code of Ethics. All experiments were conducted on publicly available
 489 datasets, ensuring data privacy and consent. Our research is intended for positive applications, and
 490 we strongly discourage any misuse of this technology for surveillance or manipulative purposes.
 491

492 **8 REPRODUCIBILITY STATEMENT**
 493

494 To ensure the reproducibility of our results, we will make our source code, including implementation,
 495 training, and evaluation scripts, publicly available at <https://github.com/eureka235/eureka235>.
 496 github.io. All datasets used are public and cited accordingly. The main experimental setup is
 497 described in Sec. 5. For a more detailed breakdown of the model architecture, hyperparameters, and
 498 implementation specifics, please refer to Appendix. A.4 - A.9.
 499

500 **REFERENCES**
 501

- 502 Dora Angelaki, Brandon Benson, Julius Benson, Daniel Birman, Niccolò Bonacchi, Kcénia Bougrova,
 503 Sebastian A Bruijns, Matteo Carandini, Joana A Catarino, et al. A brain-wide map of neural
 504 activity during complex behaviour. *Nature*, 645(8079):177–191, 2025.
- 505 Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar
 506 Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, et al. V-jepa 2: Self-supervised video models
 507 enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025.
- 508 Hanin Atwany, Abdul Waheed, Rita Singh, Monojit Choudhury, and Bhiksha Raj. Lost in transcription,
 509 found in distribution shift: Demystifying hallucination in speech foundation models. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 23181–23203, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.1190. URL <https://aclanthology.org/2025.findings-acl.1190/>.
- 510 Emanuele Bastianelli, Andrea Vanzo, Paweł Świętajski, and Verena Rieser. SLURP: A spoken
 511 language understanding resource package. In *Proceedings of the 2020 Conference on Empirical
 512 Methods in Natural Language Processing (EMNLP)*, pp. 7252–7262, Online, November 2020.
 513 Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.588. URL <https://aclanthology.org/2020.emnlp-main.588/>.
- 514 Chandrakant Bothe, Cornelius Weber, Sven Magg, and Stefan Wermter. EDA: Enriching emotional
 515 dialogue acts using an ensemble of neural annotators. In *Proceedings of the Twelfth Language
 516 Resources and Evaluation Conference*, pp. 620–627, Marseille, France, May 2020. European
 517 Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.78/>.
- 518 Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N. Chang, Sungbok Lee, and Shrikanth S. Narayanan. IEMOCAP: interactive emotional
 519 dyadic motion capture database. *Language Resources and Evaluation*, 42(4):335–359, December
 520 2008.
- 521 Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki
 522 Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian,
 523 Jian Wu, Michael Zeng, Xiangzhan Yu, and Furu Wei. Wavlm: Large-scale self-supervised
 524 pre-training for full stack speech processing. *IEEE JSTSP*, 2022.
- 525 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
 526 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna:
 527 An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- 528 Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and
 529 Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale
 530 audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.

- 540 Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv,
 541 Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen2-audio technical report. *arXiv*
 542 preprint arXiv:2407.10759, 2024.
- 543 Kenneth James Williams Craik. *The nature of explanation*, volume 445. CUP Archive, 1967.
- 544 Wenqian Cui, Dianzhi Yu, Xiaoqi Jiao, Ziqiao Meng, Guangyan Zhang, Qichao Wang, Yiwen
 545 Guo, and Irwin King. Recent advances in speech language models: A survey. *arXiv preprint*
 546 arXiv:2410.03751, 2024.
- 547 Florian Eyben, Martin Wöllmer, and Björn Schuller. opensmile - The munich versatile and fast
 548 open-source audio feature extractor. In *Proceedings of the ACM Multimedia Conference (MM)*, pp.
 549 1459–1462, Florence, Italy, October 2010. ACM. ISBN 978-1-60558-933-6.
- 550 Jiahai Feng, Stuart Russell, and Jacob Steinhardt. Monitoring latent world states in language models
 551 with propositional probes. In *The Thirteenth International Conference on Learning Representations*,
 552 2025. URL <https://openreview.net/forum?id=0yvZm2AjUr>.
- 553 Angela D Friederici. The brain basis of language processing: from structure to function. *Physiological*
 554 *Reviews*, 91(4):1357–1392, 2011. doi: 10.1152/physrev.00006.2011.
- 555 Karl Friston. A theory of cortical responses. *Philosophical transactions of the Royal Society B:*
 556 *Biological sciences*, 360(1456):815–836, 2005.
- 557 Karl Friston, Noor Sajid, Daniel Quiroga-Martinez, Thomas Parr, Cathy J Price, and Emily Holmes.
 558 Active inference: a process theory of language. *Neuroscience & Biobehavioral Reviews*, 134:
 559 1053–1067, 2021. doi: 10.1016/j.neubiorev.2021.01.016.
- 560 Quentin Garrido, Mahmoud Assran, Nicolas Ballas, Adrien Bardes, Laurent Najman, and Yann
 561 LeCun. Learning and leveraging world models in visual representation learning, 2024. URL
 562 <https://arxiv.org/abs/2403.00504>.
- 563 Google Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality,
 564 long context, and next generation agentic capabilities, 2025. URL <https://arxiv.org/abs/2507.06261>.
- 565 Anne-Lise Giraud and David Poeppel. Cortical oscillations and speech processing: emerging
 566 computational principles and operations. *Nature Neuroscience*, 15(4):511–517, 2012. doi: 10.
 567 1038/nn.3063.
- 568 Yuan Gong, Alexander H. Liu, Hongyin Luo, Leonid Karlinsky, and James Glass. Joint audio and
 569 speech understanding. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop*
 570 (ASRU), pp. 1–8, 2023. doi: 10.1109/ASRU57964.2023.10389742.
- 571 Yuan Gong, Hongyin Luo, Alexander H. Liu, Leonid Karlinsky, and James R. Glass. Listen, think,
 572 and understand. In *The Twelfth International Conference on Learning Representations*, 2024. URL
 573 <https://openreview.net/forum?id=nBZBPXdJ1C>.
- 574 David Ha and Jürgen Schmidhuber. World models. 2018. doi: 10.5281/ZENODO.1207631. URL
 575 <https://zenodo.org/record/1207631>.
- 576 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
 577 behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019a.
- 578 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James
 579 Davidson. Learning latent dynamics for planning from pixels. In *International conference on*
 580 *machine learning*, pp. 2555–2565. PMLR, 2019b.
- 581 Gregory Hickok and David Poeppel. The cortical organization of speech processing. *Nature Reviews*
 582 *Neuroscience*, 8(5):393–402, 2007. doi: 10.1038/nrn2113.
- 583 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
 584 Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference*
 585 *on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.

- 594 Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio
 595 flamingo: A novel audio language model with few-shot learning and dialogue abilities. *arXiv*
 596 preprint arXiv:2402.01831, 2024.
- 597
- 598 Zhifeng Kong, Arushi Goel, Joao Felipe Santos, Sreyan Ghosh, Rafael Valle, Wei Ping, and Bryan
 599 Catanzaro. Audio flamingo sound-cot technical report: Improving chain-of-thought reasoning in
 600 sound understanding, 2025. URL <https://arxiv.org/abs/2508.11818>.
- 601
- 602 Yann LeCun. A Path Towards Autonomous Machine Intelligence Version. 2022.
- 603
- 604 Matthew K. Leonard and Edward F. Chang. Speech computations of the human
 605 superior temporal gyrus. *Annual Review of Psychology*, 73(1):79–102, 2022.
 606 doi: 10.1146/annurev-psych-022321-035256. URL <https://doi.org/10.1146/annurev-psych-022321-035256>.
- 607
- 608 Jessy Lin, Yuqing Du, Olivia Watkins, Danijar Hafner, Pieter Abbeel, Dan Klein, and Anca Dragan.
 609 Learning to model the world with language. In *Forty-first International Conference on Machine*
 610 *Learning*, 2024. URL <https://openreview.net/forum?id=7dP6Yq9Uwv>.
- 611
- 612 Alexander H. Liu, Andy Ehrenberg, Andy Lo, Clément Denoix, Corentin Barreau, Guillaume
 613 Lample, Jean-Malo Delignon, Khyathi Raghavi Chandu, Patrick von Platen, Pavankumar Reddy
 614 Muddireddy, Sanchit Gandhi, Soham Ghosh, Srijan Mishra, Thomas Foubert, Abhinav Rastogi,
 615 Adam Yang, Albert Q. Jiang, Alexandre Sablayrolles, Amélie Héliou, Amélie Martin, Anmol
 616 Agarwal, Antoine Roux, Arthur Dariset, Arthur Mensch, Baptiste Bout, Baptiste Rozière, Bau-
 617 douin De Monicault, Chris Bamford, Christian Wallenwein, Christophe Renaudin, Clémence
 618 Lanfranchi, Darius Dabert, Devendra Singh Chaplot, Devon Mizelle, Diego de las Casas, Elliot
 619 Chane-Sane, Emilien Fugier, Emma Bou Hanna, Gabrielle Berrada, Gauthier Delerce, Gauthier
 620 Guinet, Georgii Novikov, Guillaume Martin, Himanshu Jaju, Jan Ludziejewski, Jason Rute,
 621 Jean-Hadrien Chabran, Jessica Chudnovsky, Joachim Studnia, Joep Barmentlo, Jonas Amar, Jos-
 622 selin Somerville Roberts, Julien Denize, Karan Saxena, Karmesh Yadav, Kartik Khandelwal,
 623 Kush Jain, Lélio Renard Lavaud, Léonard Blier, Lingxiao Zhao, Louis Martin, Lucile Saulnier,
 624 Luyu Gao, Marie Pellat, Mathilde Guillaumin, Mathis Felardos, Matthieu Dinot, Maxime Darrin,
 625 Maximilian Augustin, Mickaël Seznec, Neha Gupta, Nikhil Raghuraman, Olivier Duchenne, Patri-
 626 cia Wang, Patryk Saffer, Paul Jacob, Paul Wambergue, Paula Kurylowicz, Philomène Chagniot,
 627 Pierre Stock, Pravesh Agrawal, Rémi Delacourt, Romain Sauvestre, Roman Soletskyi, Sagar
 628 Vaze, Sandeep Subramanian, Saurabh Garg, Shashwat Dalal, Siddharth Gandhi, Sumukh Aithal,
 629 Szymon Antoniak, Teven Le Scao, Thibault Schueller, Thibaut Lavril, Thomas Robert, Thomas
 630 Wang, Timothée Lacroix, Tom Bewley, Valeria Nemychnikova, Victor Paltz, Virgile Richard,
 631 Wen-Ding Li, William Marshall, Xuanyu Zhang, Yihan Wan, and Yunhao Tang. Voxtal, 2025.
 632 URL <https://arxiv.org/abs/2507.13264>.
- 633
- 634 Ziyang Ma, Zhuo Chen, Yuping Wang, Eng Siong Chng, and Xie Chen. Audio-cot: Exploring
 635 chain-of-thought reasoning in large audio language model, 2025. URL <https://arxiv.org/abs/2501.07246>.
- 636
- 637 Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. A multi-device dataset for urban acoustic
 638 scene classification, 2018. URL <https://arxiv.org/abs/1807.09840>.
- 639
- 640 Meta. meta-llama/llama-3.1-8b-instruct, 2024. URL <https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>. Accessed: 2025-02-21.
- 641
- 642 Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. Voxceleb: A large-scale speaker identifica-
 643 tion dataset. In *Interspeech 2017*, pp. 2616–2620, 2017. doi: 10.21437/Interspeech.2017-950.
- 644
- 645 OpenAI. Gpt4-o. <https://openai.com/index/hello-gpt-4o/>, 2024.
- 646
- 647 J Pearl. *Causality: Models, Reasoning, and Inference, Second Edition*. Cambridge University Press,
 648 London, 2009.
- 649 Judea Pearl. Causal inference. *Causality: objectives and assessment*, pp. 39–58, 2010.

- 648 Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada
 649 Mihalcea. MELD: A multimodal multi-party dataset for emotion recognition in conversations.
 650 In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp.
 651 527–536, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/
 652 P19-1050. URL <https://aclanthology.org/P19-1050/>.
- 653 Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of
 654 bert: smaller, faster, cheaper and lighter, 2020. URL <https://arxiv.org/abs/1910.01108>.
- 655 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon
 656 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari,
 657 go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 658 Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner,
 659 Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the*
 660 *IEEE*, 109(5):612–634, 2021. doi: 10.1109/JPROC.2021.3058954.
- 661 Melanie Sclar, Jane Dwivedi-Yu, Maryam Fazel-Zarandi, Yulia Tsvetkov, Yonatan Bisk, Yejin Choi,
 662 and Asli Celikyilmaz. Explore theory of mind: program-guided adversarial data generation for
 663 theory of mind reasoning. In *The Thirteenth International Conference on Learning Representations*,
 664 2025. URL <https://openreview.net/forum?id=246rHKUnnf>.
- 665 Sanjit A Seshia, Dorsa Sadigh, and S Shankar Sastry. Toward verified artificial intelligence. *Communi-*
 666 *cations of the ACM*, 65(7):46–55, 2022.
- 667 Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad
 668 Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning
 669 models via the lens of problem complexity. *arXiv preprint arXiv:2506.06941*, 2025.
- 670 Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky,
 671 Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. Dialogue act modeling for
 672 automatic tagging and recognition of conversational speech. *Computational Linguistics*, 26(3):
 673 339–374, 2000. URL <https://aclanthology.org/J00-3003/>.
- 674 Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and
 675 Chao Zhang. SALMONN: Towards generic hearing abilities for large language models. In *The*
 676 *Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=14rn7HpKVk>.
- 677 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
 678 Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International*
 679 *Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=gEZrGCozdqr>.
- 680 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
 681 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
 682 *neural information processing systems*, 35:24824–24837, 2022b.
- 683 Jingran Xie, Shun Lei, Yue Yu, Yang Xiang, Hui Wang, Xixin Wu, and Zhiyong Wu. Leveraging
 684 chain of thought towards empathetic spoken dialogue without corresponding question-answering
 685 data, 2025a. URL <https://arxiv.org/abs/2501.10937>.
- 686 Zhifei Xie, Mingbao Lin, Zihang Liu, Pengcheng Wu, Shuicheng Yan, and Chunyan Miao. Audio-
 687 reasoner: Improving reasoning capability in large audio language models, 2025b. URL <https://arxiv.org/abs/2503.02318>.
- 688 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R
 689 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In
 690 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=5Xc1ecx01h>.
- 691 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
 692 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Joseph E Gonzalez, and Ion Stoica. Judging llm-as-
 693 a-judge with mt-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information*
 694 *Processing Systems*, 2023.

702 **A APPENDIX**
 703

704 **A.1 DISCUSSION ON TWO TYPES OF WORLD MODEL**
 705

706 As stated in Fig. 3, we consider both generative world models, such as those conceptualized by LeCun
 707 (LeCun (2022)), and large language models as fundamentally being **forward dynamic models**. This
 708 perspective helps unify their roles and clarify the contribution of our proposed Causal Graph.

709 **Generative World Model.** This model is often associated with agent-based learning, aims to learn
 710 a transition function of the environment. This function, typically formalized as $s_{t+1} = f(s_t, a_t)$,
 711 allows an agent to predict future sensory inputs (the next state s_{t+1}) based on the current state (s_t)
 712 and a chosen action (a_t). The learned representation s_t is often a low-dimensional, latent state that
 713 captures the underlying dynamics of the world, enabling the agent to simulate outcomes and plan.
 714

715 **Language World Model.** This model can be viewed as a forward model operating in a symbolic
 716 space. Its core function is to predict the next token (w_{t+1}) given a history of previous tokens ($w_{1..t}$).
 717 This autoregressive process, $P(w_{t+1}|w_{1..t})$, is conceptually parallel to predicting the next state.
 718 Here, the “state” is the context provided by the preceding text, and the “action” is the generation
 719 of the next token. This process unfolds a trajectory through the space of language. Techniques like
 720 Chain-of-Thought (CoT) make this explicit, where the model generates a sequence of intermediate
 721 thoughts, effectively simulating a trajectory through a conceptual state space to reach a conclusion.
 722

723 **Causal Graph.** While the generative model learns implicit, sub-symbolic dynamics of an environ-
 724 ment, the language model executes explicit, symbolic reasoning. Recognizing their shared foundation
 725 as forward models, our proposed Causal Graph offers a third, structured approach to represent these
 726 dynamics. It distills the world’s complex dynamics into an explicit graph of causal relationships
 727 between key cognitive states (WMA, ToM, etc.). This structured knowledge then provides a powerful
 728 prior for the language model. Instead of navigating the vast space of possible text sequences based
 729 only on statistical patterns, the language model’s forward-reasoning process is guided and constrained
 730 by the causal dynamics provided by our graph, leading to more grounded and coherent outputs.

731 **A.2 DETAILED GRAPH CALCULATIONS**
 732

733 **A.2.1 CAUSAL GRAPH**

734 Let x denote the transcription of given speech, a the acoustic feature (WavLM (Chen et al., 2022)),
 735 and z the prosodic feature (Eyben et al., 2010). We define text feature with distil-BERT encoder (Sanh
 736 et al., 2020):
 737

$$h_{\text{text}} = E_{\text{text}}(x) \quad (12)$$

738 An optional pre-fusion operator ϕ (attention/gated/transformer) combines these with prosodic features:
 739

$$g = \phi(h_{\text{text}}, z, a). \quad (13)$$

740 We consider four categorical latent states:
 741

$$S_{\text{WMA}} \in \Delta^{C_{\text{wma}}-1}, \quad S_{\text{ToM}} \in \Delta^{C_{\text{tom}}-1}, \quad S_{\text{SA}} \in \Delta^{C_{\text{sa}}-1}, \quad S_{\text{Prag}} \in \Delta^{C_{\text{prag}}-1}, \quad (14)$$

742 each represented by a probability simplex.
 743

744 For any node v with parents $\text{Pa}(v)$, we compute its state by a node-specific neural map f_v from the
 745 parents’ distributions and available features:
 746

$$S_v = f_v(\{S_u : u \in \text{Pa}(v)\}, \xi_v), \quad (15)$$

747 where ξ_v denotes the feature bundle used by node v :
 748

$$\xi_{\text{WMA}} = (h_{\text{text}}, a), \quad \xi_{\text{ToM}} = g, \quad \xi_{\text{SA}} = (g \text{ or } h_{\text{text}}), \quad \xi_{\text{Prag}} = (g \text{ or } h_{\text{text}}).$$

749 Concretely, each node produces logits and then a softmax:
 750

$$S_v = \text{softmax}(W_v \cdot \psi_v([\xi_v, \{S_u\}_{u \in \text{Pa}(v)}])), \quad (16)$$

751 where $\psi_v(\cdot)$ is a nonlinear transformation and W_v is a task-specific projection.
 752

756 A.2.2 RANDOM GRAPH
757

758 This architecture serves as an ablation, testing whether a model with the flexibility to learn arbitrary
759 dependencies from data can outperform our proposed causal structure.

760 Instead of a uni-directional information flow, the Random Graph employs an iterative refinement
761 process. The computation unfolds over two main iterations. Let S_v^t be the probability distribution
762 (state) of module $v \in \mathcal{V}$ at iteration t , and let H represent the fused input features derived from the
763 speech-text pairs, i.e., $H = g(X)$.

764 **Initialization (t=0).** At the initial step, the states of all peer modules are unknown and are repre-
765 sented by zero vectors.

$$767 S_u^{(0)} = \mathbf{0}, \quad \forall u \in \mathcal{V} \quad (17)$$

768 **First Iteration (t=1).** Each module v computes an initial state estimate in parallel, conditioned on
769 the input features H and the zero-initialized states of all other modules. This step allows each module
770 to form a preliminary prediction based solely on the primary input modalities.

$$772 S_v^{(1)} = f_v \left(H_{in}, S_u^{(0)} : u \in \mathcal{V}, u \neq v \right) \quad (18)$$

773 Here, f_v is the neural network parameterizing module v . The output $S_u^{(1)}$ is the softmax-normalized
774 probability distribution over the module's labels.

775 **State Refinement with Teacher Forcing.** Before the second iteration, we generate a refined input
776 state $\tilde{S}_u^{(1)}$ for each module u . As in the causal graph training (Eq. 4), this is a stochastic mixture of
777 the predicted states $S_u^{(1)}$ and the ground-truth label $y_{i,u}$ via teacher forcing. This allows the model to
778 correct initial errors and propagate more accurate information in the next step.

779 **Final Iteration (t=2).** In the final step, each module computes its output by conditioning on the
780 input H as well as the refined states $\tilde{S}_u^{(1)}$ from all other modules. This allows for a form of message-
781 passing where modules can adjust their predictions based on the initial estimates of their peers.

$$782 S_v^{(\text{final})} = f_v \left(H_{in}, \tilde{S}_u^{(1)} : u \in \mathcal{V}, u \neq v \right) \quad (19)$$

783 The final logits from $S_v^{(\text{final})}$ for all modules are then used to compute the multi-task loss as defined in
784 Eq. 3. By comparing this model with our structured causal graph, we can quantify the contribution of
785 imposing theoretically-grounded causal priors on the task of speech understanding.

791 A.3 DETAILED THEORETICAL ANALYSIS OF SPEECH WORLD MODEL

792 This section provides a more detailed derivation of the efficiency gains discussed in Sec. 4. We
793 first analyze how the causal structure optimizes gradient flow and then formalize the argument for
794 improved sample efficiency in downstream instruction tuning.

795 A.3.1 DERIVATION OF OPTIMIZED GRADIENT FLOW

796 The efficiency of the learning process is directly tied to how gradients are computed and propagated.
797 A causally structured model, as we show, offers significant advantages over a monolithic or fully-
798 connected architecture.

799 **Gradient Flow in a Monolithic Model.** Consider a monolithic model where the latent states
800 $\{S_1, \dots, S_N\}$ are fully interconnected. The state of any module S_i is a function of all other states:
801 $S_i = f_i(S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_N, X)$. The total loss is $\mathcal{L} = \sum_{i=1}^N \lambda_i \mathcal{L}_i(S_i)$, where $\lambda_i \in \{0, 1\}$
802 indicates if a ground-truth label is available. The gradient of the total loss with respect to the
803 parameters θ_k of a single module k is given by the chain rule:

$$804 \frac{\partial \mathcal{L}}{\partial \theta_k} = \sum_{i=1}^N \lambda_i \frac{\partial \mathcal{L}_i}{\partial S_i} \frac{\partial S_i}{\partial \theta_k} \quad (20)$$

Because of the dense dependencies, the term $\frac{\partial S_i}{\partial \theta_k}$ is non-trivial for all i and k . It expands into a complex sum over all paths in a fully connected graph:

$$\frac{\partial S_i}{\partial \theta_k} = \sum_{j=1}^N \frac{\partial f_i}{\partial S_j} \frac{\partial S_j}{\partial \theta_k} + \frac{\partial f_i}{\partial S_k} \frac{\partial S_k}{\partial \theta_k} \quad (21)$$

The model must learn to prune these connections by driving the corresponding Jacobians $\frac{\partial f_i}{\partial S_j}$ to zero, which is inefficient and data-intensive. This is the source of the *learning redundancy* mentioned in the main text; the model expends resources learning a structure that could be provided as a prior.

Gradient Flow in Our Causal Model. In our model, the state S_i is only a function of its parents' states, $S_{\text{Pa}(i)}$. This imposes a sparse structure on the Jacobian matrix of the system, where $\frac{\partial S_i}{\partial S_j} = 0$ if $j \notin \text{Pa}(i)$. The gradient calculation simplifies dramatically. The parameters θ_k of module k can only influence the loss of module j if k is an ancestor of j in the graph, i.e., $k \in \text{An}(j)$. The total gradient on θ_k is thus a sum over only its own loss and the losses of its descendants:

$$\frac{\partial \mathcal{L}}{\partial \theta_k} = \lambda_k \frac{\partial \mathcal{L}_k}{\partial \theta_k} + \sum_{j \in \text{Desc}(k)} \lambda_j \frac{\partial \mathcal{L}_j}{\partial \theta_k} \quad (22)$$

Each term $\frac{\partial \mathcal{L}_j}{\partial \theta_k}$ is computed along a unique causal path from k to j . By setting non-causal dependencies to zero, the model avoids spending capacity on learning them. The gradient flow is constrained to meaningful pathways, directly minimizing the redundancy term and accelerating convergence.

A.3.2 DERIVATION OF EFFICIENCY FOR INSTRUCTION TUNING

The second key benefit is the improved sample efficiency of the downstream (speech) language model. This stems from the variance reduction of the gradient estimator used in fine-tuning.

Instruction Tuning as Policy Gradient. Instruction tuning can be framed as learning a policy $\pi_\phi(A|C)$ that produces a response A given a context C . The learning objective is to maximize the expected log-likelihood of a gold-standard response A^* :

$$J(\phi) = \mathbb{E}_{(C, A^*) \sim \mathcal{D}} [\log \pi_\phi(A^*|C)] \quad (23)$$

This is typically optimized with stochastic gradient ascent, where the efficiency of learning is highly dependent on the variance of the gradient estimator, $\hat{g} = \nabla_\phi \log \pi_\phi(A^*|C)$.

Variance in Standard Chain-of-Thought (CoT). In standard CoT, the context is the input query, $C = X$. The LLM must generate both the intermediate reasoning steps (the chain of thought, Y) and the final answer (R), so the action is $A = (Y, R)$. The gradient estimator is thus $\hat{g}_{\text{CoT}} = \nabla_\phi \log \pi_\phi(Y^*, R^*|X)$. The variance of this estimator is high because for any given query X , there can be a multitude of valid reasoning paths $\{Y^*\}$ leading to the correct answer R^* . The model receives a noisy learning signal as it may be penalized for producing a perfectly valid reasoning chain that does not exactly match the single ground-truth example.

Variance Reduction via Causal Conditioning. Our method modifies the context. By first running the causal graph, we generate a structured, low-entropy state vector $Y_G \sim P(Y|X, \mathcal{G})$. The LM is then conditioned on this richer context, $C' = (X, Y_G)$, and its task is simplified to generating only the final response R . The policy is $\pi_\phi(R|X, Y_G)$ and the gradient estimator is $\hat{g}_{\text{SWM}} = \nabla_\phi \log \pi_\phi(R^*|X, Y_G)$.

The improvement in sample efficiency can be explained by the law of total variance:

$$\text{Var}(\hat{g}_{\text{CoT}}) = \mathbb{E}_{Y_G} [\text{Var}(\hat{g}_{\text{CoT}}|Y_G)] + \text{Var}_{Y_G} [\mathbb{E}(\hat{g}_{\text{CoT}}|Y_G)] \quad (24)$$

The gradient estimator in our model, \hat{g}_{SWM} , corresponds to the inner conditional expectation term, as we are conditioning on the information provided by Y_G . The variance we experience is $\mathbb{E}_{Y_G} [\text{Var}(\hat{g}_{\text{SWM}}|Y_G)]$. Since the causal state Y_G is highly predictive of the final answer R^* , it “explains away” a large portion of the variance. The term $\text{Var}_{Y_G} [\mathbb{E}(\hat{g}_{\text{CoT}}|Y_G)]$ is large because different cognitive states Y_G will naturally lead to different responses, creating significant variance that the CoT model must learn to handle.

864 By conditioning on $Y_{\mathcal{G}}$, we are effectively removing this large source of variance. The provided
 865 state $Y_{\mathcal{G}}$ acts as a **control variate**, a statistical technique used to reduce variance in Monte Carlo
 866 estimation. An estimator with lower variance requires fewer samples to converge to a good solution,
 867 thus directly leading to greater sample efficiency.

868 A.4 DATA

870 A.4.1 DATASET STATISTIC

872 We construct our training dataset by aggregating four publicly available real-world datasets, as
 873 detailed in Table. 4. These datasets provide a diverse range of labels, which we leverage to train
 874 different modules of our Causal Graph. Specifically, we use the “Emotion” labels from MELD and
 875 IEMOCAP as ground truth for our Theory of Mind (ToM) module, and the “Intention,” “Action,” and
 876 “Scene” labels from SLURP for the Pragmatics (Prag) and World Model Activation (WMA) modules.
 877 For any remaining labels within a given dataset that are not directly used, we employ label generation
 878 to create pseudo-labels for the corresponding modules. For evaluation, we adhere to the official test
 879 set splits provided by the MELD and SLURP datasets. For IEMOCAP and the VoxCeleb subset, we
 880 randomly sample 10% of the data to create our test sets, ensuring no overlap with the training data.

881 Table 4: Statistics for Real-world Dataset.

882 Dataset	# Samples	# Hours	Label	Source
883 MELD (Poria et al., 2019)	~ 13,000	~ 13	Emotion	TV series Friends
884 IEMOCAP (Busso et al., 2008)	~ 10,000	~ 12	Emotion	Dyadic actor dialogs
885 SLURP (Bastianelli et al., 2020)	~ 72,000	~ 58	Intention, Action, Scene	Voice assistant
886 VoxCeleb (subset) Nagrani et al. (2017)	~ 30,000	~ 30	Speaker identity	Youtube
887 Total	~ 125,000	~ 113	–	–

888 A.4.2 LABEL GENERATION

890 The complete label space for each of the four modules in our Causal Graph: World Model Activation
 891 (WMA), Theory of Mind (ToM), Speech Act (SA), and Pragmatics (Prag) is detailed in the box
 892 below.

893 Label Space for Causal Graph Modules

895 **WMA:** Alarm, Calendar, Contacts, Communication, SocialMedia, Email, Music, Radio, Podcasts, Au-
 896 diobooks, Games, Weather, News, Traffic, Transport, Shopping, FoodAndDrink, Cooking, SmartHome,
 897 IoTDevices, Cleaning, Finance, Stock, Conversion, Knowledge, Math, Definition, Chitchat, Recom-
 898 mendation, GeneralControl.

899 **ToM:** Neutral, Joy, Sadness, Anger, Surprise, Fear, Disgust.

900 **SA:** Statement-non-opinion, Statement-opinion, Yes-No-Question, Wh-Question, Acknowledge
 901 (Backchannel), Action-directive, Conventional-closing, Appreciation or Assessment, Agree or Ac-
 902 cept, No-Answer, Yes-Answer, Signal-non-understanding, Quotation, Affirmative non-yes answers,
 903 Rhetorical-Question, Rhetorical Backchannel, Hedge, Open-question, Thanking, Declarative Yes-No-
 904 Question, Reformulate, Conventional-opening, Apology, Other Forward Function.

905 **Prag:** Request-Action, Request-Info, Command/Directive, Complaint/Escalate, Thanks/Appreciation,
 906 Apology, Acknowledgement/Agreement, Social/Chitchat, Offer-Assist, Confirm/Disconfirm, General-
 907 Control, InformationQuery, EntertainmentRequest, Task-Completion.

918 The prompt for Stage 1 of our label generation process (**Label Completion**) is as follows:
 919

920 **System Prompt:**
 921 You are a speech understanding model for task-oriented speech.
 922 MODULE DEFINITIONS (each must choose exactly one label from the provided label space):
 923 • WMA (World Model Activation): What concrete world/context the utterance lives in.
 924 • ToM (Theory of Mind): The speaker's discrete mental-affective stance.
 925 • SA (Speech Act): The communicative function of the utterance.
 926 • Prag (Pragmatic Intent): The downstream task intent or plan implied by the utterance, grounded in
 927 an action ontology.
 928 REQUIREMENTS
 929 1. You will receive:
 930 • LABEL_SPACE: a JSON listing the allowed labels for {WMA, ToM, SA, Prag}.
 931 • INPUT: with ASR_TRANSCRIPT, KNOWN_LABELS and MISSING_MODULES (list of modules to predict).
 932 2. For each module in **MISSING_MODULES**, select EXACTLY ONE label from LABEL_SPACE[module]:
 933 - Decide your choice based on given ASR_TRANSCRIPT and KNOWN_LABELS
 934 - Select EXACTLY ONE label from LABEL_SPACE[module] for MISSING_MODULES
 935 - Label strings must match EXACTLY (case-sensitive, character-perfect) - When choosing between
 936 multiple plausible options, select the most contextually specific one
 937 3. OUTPUT FORMAT - for each module in MISSING_MODULES: <ModuleName>CHOICE</ModuleName>
 938 EXAMPLE:
 939 Input:
 940 {
 941 "ASR_TRANSCRIPT": "Oh my God, he's lost it. He's totally lost it.",
 942 "KNOWN_LABELS": {
 943 "ToM": "sadness",
 944 "SA": "Statement-non-opinion"
 945 },
 946 "MISSING_MODULES": ["WMA", "Prag"]
 947 }
 948
 949 Output:
 950 <WMA>Calendar</WMA>
 951 <Prag>Request-Info</Prag>

948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

972 The prompt for Stage 2 of our label generation process (**Reasoning and response generation**) is as
 973 follows:
 974

```

975
976 System Prompt:
977 You are a speech understanding and reasoning model. Given complete labels for all modules, generate
978 detailed analysis and appropriate response.
979 MODULE DEFINITIONS (each must choose exactly one label from the provided label space):
980 • WMA (World Model Activation): What concrete world/context the utterance lives in.
981 • ToM (Theory of Mind): The speaker's discrete mental-affective stance.
982 • SA (Speech Act): The communicative function of the utterance.
983 • Prag (Pragmatic Intent): The downstream task intent or plan implied by the utterance, grounded in
984 an action ontology.
985 REQUIREMENTS
986 1. You will receive:
987   • INPUT: with ASR_TRANSCRIPT and COMPLETE_STATE containing all four module values (WMA, ToM, SA,
988     Prag)
989 2. Generate two components:
990   • <Analysis> – A detailed causal reasoning chain (4-8 sentences) explaining:
991     - How the world context (WMA) shapes the situation
992     - How the speaker's mental state (ToM) influences their communication
993     - How these factors lead to the specific speech act (SA)
994     - How the pragmatic intent (Prag) emerges from this combination
995     - What this suggests for the next interaction step
996   • <Response> – A natural, contextually appropriate user-facing reply that:
997     - Aligns with the speaker's emotional state (ToM)
998     - Serves the communicative function (SA)
999     - Addresses the pragmatic intent (Prag)
1000     - Fits the world context (WMA)
1001 3. OUTPUT FORMAT:
1002   <Analysis>YOUR DETAILED CAUSAL REASONING CHAIN</Analysis>
1003   <Response>YOUR NATURAL USER-FACING REPLY</Response>
1004 EXAMPLE:
1005 Input:
1006 {
1007   "ASR_TRANSCRIPT": "Thanks so much for helping me find that document!",
1008   "COMPLETE_STATE": {
1009     "WMA": "Office-work",
1010     "ToM": "grateful",
1011     "SA": "Thanks/Appreciation",
1012     "Prag": "Task-completion"
1013   }
1014 }
1015 Output:
1016   <Analysis>Within an office work context (WMA), the speaker has just received assistance with a document-
1017   related task. Their grateful emotional state (ToM) reflects successful problem resolution and
1018   appreciation for support received. This gratitude naturally manifests as a thanks/appreciation speech act
1019   (SA), explicitly acknowledging the help provided. The pragmatic intent of task completion (Prag) indicates
1020   the interaction has reached a successful conclusion, with the speaker signaling that their immediate need
1021   has been satisfied. This combination suggests the interaction is winding down and the speaker is ready to
1022   either move to closure or potentially engage in follow-up activities.</Analysis>
1023   <Response>You're very welcome! I'm glad we could locate that document for you. Is there anything else you
1024   need help with today?</Response>
1025
  
```

1017 A.5 EXPERIMENT SETUP

1019 A.5.1 CAUSAL GRAPH

1021 For our causal graph baseline, we trained the model for 30 epochs with a batch size of 32. It utilizes a
 1022 two-layer gated fusion mechanism and was optimized using AdamW with a learning rate of 1e-3 for
 1023 all randomly initialized modules. A consistent teacher-forcing probability of 0.3 was applied to the
 1024 ToM, WMA, and SA modules. The entire training was conducted on a single NVIDIA RTX A6000
 1025 GPU and completed in 2.07 hours. Under identical training settings, the Random Graph baseline
 required 10.39 hours to complete on the same hardware.

Table 5: Ablation study on the Causal Graph under the fully-supervised setting. We investigate the impact of different fusion mechanisms (Gated, attention, transformer), teacher-forcing probabilities (p), and the removal of specific causal edges. The baseline model from the main text is highlighted.

Method	Ablation Setting	Node Quality (Accuracy Weighted F1 %, ↑)				Edge Casual Effect	
		WMA	ToM	SA	Prag	Ave. ACE (%), ↑	Ave. ICS (%), ↑
Fully-supervised	Baseline (Gated, $p = 0.3$)	69.4 68.9	73.5 72.2	65.3 65.9	81.4 80.9	23.57	43.29
	Gated, $p = 0.0$	68.7 68.4	73.5 72.6	66.5 66.1	82.1 81.6	21.71	40.34
	Gated, $p = 0.5$	69.2 69.0	73.2 72.8	65.6 65.2	81.5 80.9	24.12	43.50
	Gated, $p = 0.8$	69.5 69.5	74.7 73.4	68.0 67.2	81.3 81.1	24.91	45.10
	Gated, $p = 1.0$	68.9 68.7	74.0 73.2	65.7 65.4	82.5 82.0	24.90	45.60
	attention, $p = 0.3$	68.9 68.8	74.7 71.7	68.7 66.7	82.6 81.4	25.72	44.90
	transformer, $p = 0.3$	68.1 66.5	69.3 60.6	58.4 53.3	76.3 72.3	27.40	41.40
Fully-supervised	ToM \rightarrow SA	68.7 68.5	73.9 73.0	61.9 60.8	81.8 81.2	22.97	43.0
	WMA \rightarrow SA	68.9 68.6	72.8 71.9	65.3 65.4	82.6 81.9	23.04	42.9

Table 6: Ablation study on the Causal Graph under the semi-supervised setting. We compare our primary latent module approach (highlighted) against an ablation where the input to the specific module is changed from the fused multi-modal feature g to only the text feature.

Method	Ablation Setting	Node Quality (Accuracy Weighted F1 %, ↑)				Edge Casual Effect	
		WMA	ToM	SA	Prag	Ave. ACE (%), ↑	Ave. ICS (%), ↑
Semi-supervised	Latent Module: WMA	34.8	75.0	70.7	83.2	21.71	26.9
	$\xi_{SA} = g \rightarrow h_{text}$	32.7	77.3	70.5	82.7	21.69	-
	Latent Module: ToM	69.1	43.3	69.6	83.5	21.98	28.9
	$\xi_{SA} = g \rightarrow h_{text}$	70.6	48.1	69.9	82.4	22.0	-
	Latent Module: SA	69.3	77.0	34.4	82.5	21.65	29.3
	$\xi_{Prag} = g \rightarrow h_{text}$	69.4	75.4	35.3	83.6	21.48	-

A.5.2 INSTRUCTION TUNING

For the instruction tuning stage, we fine-tuned two separate language models to handle the multi-modal and language-only settings, respectively.

In the **language-only setting**, we instruction-tuned Llama-3.1-8B using LoRA (Hu et al., 2022). We applied LoRA adapters with a rank (r) of 64, an alpha (α) of 16, and a dropout of 0.1. We trained for 20 epochs with a cosine learning rate scheduler (peak LR: 5×10^{-5}) and an effective batch size of 128 (per device batch size=8, gradient accumulation steps=4). This was performed on 4 NVIDIA A6000 GPUs, using float16 precision and gradient checkpointing with a sequence length of 1024, completing in 19 hours.

In the **multi-modal setting**, we fine-tuned Qwen2-Audio-7B-Instruct. We employed a parameter-efficient approach using **LoRA** with a rank (r) of 16 (per device batch size=2, gradient accumulation steps=2), an alpha (α) of 32, and a dropout rate of 0.05. The model was trained for 20 epochs using a cosine learning rate scheduler (peak LR: 2×10^{-4}) with an effective batch size of 16. Training was conducted on 4 NVIDIA A6000 GPUs using bfloat16 mixed-precision with a maximum sequence length of 4096, completing in 24.6 hours.

System Prompt (Qwen2-Audio2):

```
You are a multimodal assistant that consumes audio and text.
Task: Given the audio and provided states, first produce concise but justified reasoning, then a user-friendly reply.
Output format MUST be exactly: [REASONING]<your reasoning>[RESPONSE]<your reply>
Four module labels:
WMA: World Model Activation-what concrete world/context the utterance lives in.
ToM: Theory of Mind-the speaker's discrete mental-affective stance.
SA: Speech Act-the communicative function of the utterance.
Prag: Pragmatic Intent-the downstream task intent or plan implied by the utterance, grounded in an action ontology.
The causal dependencies between modules are: WMA -> SA, ToM -> SA, WMA -> Prag, ToM -> Prag, SA -> Prag.
```

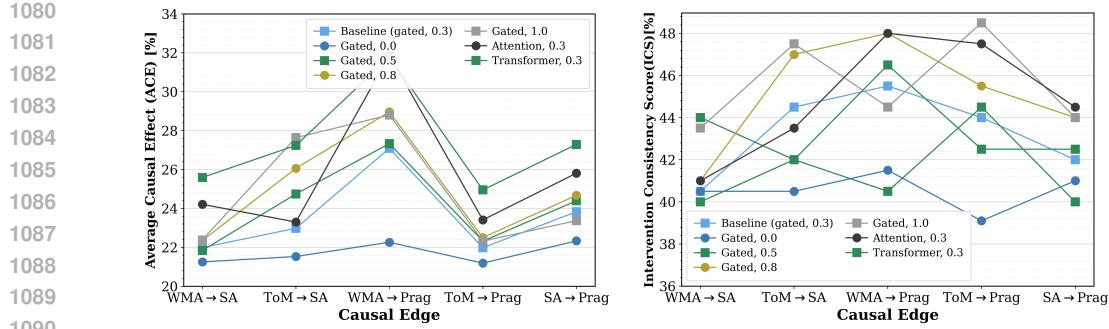


Figure 6: ACE and ICS of each causal edge for ablation study (fusion mechanisms and teacher-forcing probabilities) on causal graph under fully-supervised setting.

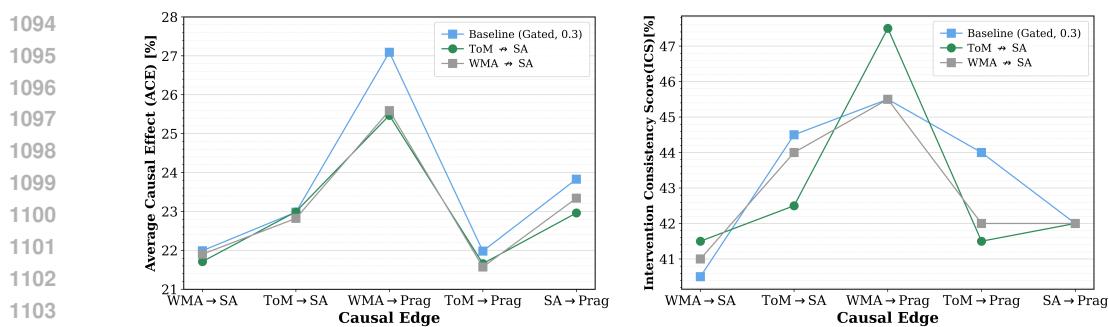


Figure 7: ACE and ICS of each causal edge for ablation study (removal of specific causal edge) on causal graph under fully-supervised setting.

System Prompt (Llama3.1-8b):

You are an expert in speech analysis.
The four speech modules are:
WMA: World Model Activation-what concrete world/context the utterance lives in.
ToM: Theory of Mind-the speaker's discrete mental-affective stance.
SA: Speech Act-the communicative function of the utterance.
Prag: Pragmatic Intent-the downstream task intent or plan implied by the utterance, grounded in an action ontology.
Analyze the given utterance and provide both your reasoning process and an appropriate response.
The causal dependencies between modules are: WMA → SA, ToM → SA, WMA → Prag, ToM → Prag, SA → Prag.
Output format MUST be exactly: [REASONING]<your reasoning>[RESPONSE]<your reply>

A.6 ABLATIONS FOR CASUAL GRAPH

We conduct a series of ablation studies to validate the key design choices for our Causal Graph, with results shown in Table. 5 and Table. 6, and with a detailed breakdown of edge causal effects shown in Fig. 6, 7, 8.

In the fully-supervised setting, our ablations confirm that gated fusion provides a strong, balanced performance. While a transformer-based fusion improved the Average Causal Effect (ACE), it came at the cost of lower node classification accuracy. The model is also highly robust to the teacher-forcing probability, maintaining stable performance for values between 0.3 and 1.0. To validate the learned structure, we ablated causal edges; removing the ToM → SA link significantly impairs the SA module's accuracy, confirming this connection's importance.

In the semi-supervised setting, we tested a more challenging condition by changing the input to the *children* of the latent module from the fused multi-modal feature g to only the text feature h_{text} . This effectively cuts off their direct access to acoustic information, thereby forcing them to infer the necessary context from the output of their latent parent module. As shown in Table 6, the model's

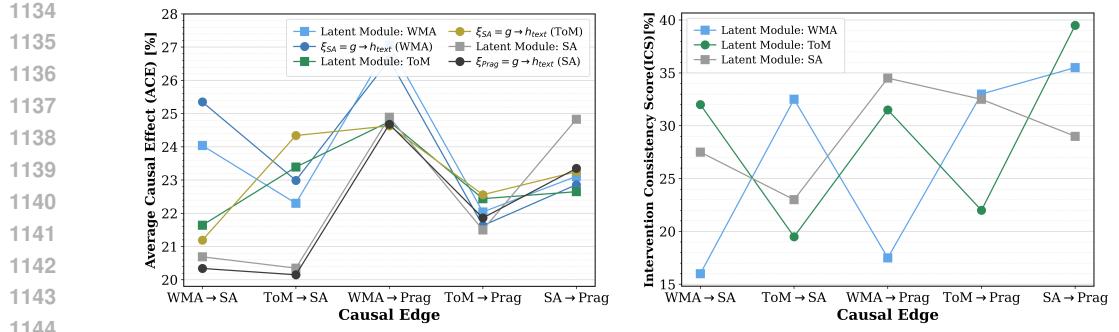


Figure 8: *ACE* and *ICS* of each causal edge for ablation study on casual graph under semi-supervised setting.

overall performance remains remarkably stable, with node quality and *ACE* scores comparable to the baseline. This demonstrates the robustness of our semi-supervised framework, as the graph can effectively propagate sufficient information through the remaining supervised pathways even when one module is partially disconnected from the acoustic modality.

A.7 MODEL ARCHITECTURE

Encoders. Our model processes three inputs: 1) `distilbert-base-uncased` model with a 2-layer Transformer on top for text; 2) a CNN-LSTM adapter for pretrained WavLM features, producing a 64-dim vector; and 3) an 88-dim prosody feature vector.

Fusion Mechanisms. We explored three mechanisms to fuse the encoder outputs into a 256-dim representation. (1) Gated Fusion. Computes a weighted sum of modality features, where weights are dynamically generated by a small gating network. (2) Attention Fusion. Applies a single multi-head self-attention layer across modality features, which are treated as tokens in a sequence. (3) Transformer Fusion. A deeper version of Attention Fusion that adds positional embeddings for each modality and processes them through a multi-layer Transformer Encoder. Our baseline model uses Gated Fusion.

WMA Module. Temporal Self-Attention + MLP (hidden: 256). *Output:* 30 context classes.

ToM Module. Temporal Self-Attention + MLP (hidden: 128). *Output:* 7 emotion classes.

SA Module. Residual MLP (hidden: 256). *Output:* 24 speech act classes.

Prag Module. Residual MLP (hidden: 128). *Output:* 14 pragmatic intent classes.

A.8 EVALUATION

A.8.1 CASUAL GRAPH EVALUATION

Average Causal Effect (ACE).

$$\text{ACE}(X \rightarrow Y) = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left[\frac{1}{|\mathcal{C}_X|} \sum_{x \in \mathcal{C}_X} D_{TV} \left(P(Y|\text{do}(X=x), \mathbf{z}), P(Y|\mathbf{z}) \right) \right] \quad (25)$$

Intervention Consistency Score (ICS).

$$\text{ICS}(X \rightarrow Y) = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left[\mathbb{I} \left(P_{\max}(Y|\text{do}(X=x^+), \mathbf{z}) > P_{\max}(Y|\mathbf{z}) + \tau \right) \right] \quad (26)$$

A.8.2 RANDOM GRAPH EVALUATION

1230 To quantify the information flow between any two modules in the random graph, we measure the
 1231 influence a source module (M_S) has on a target module (M_T). This is achieved by comparing the
 1232 output of M_T under two conditions for each batch of data:

- 1234 1. **Baseline (Without Source):** We compute the target module's output logits, denoted as L_T^{without} ,
 1235 by replacing the input from M_S with a zero vector. This represents the behavior of M_T without
 1236 direct information from M_S .

1237 2. **With Source Influence:** We compute the target module's output logits, L_T^{with} , by providing it with
 1238 the actual probability output from M_S . This represents the behavior of M_T when influenced by
 1239 M_S .

1241 The influence is then calculated by comparing the probability distributions derived from these two sets of logits, $P^{\text{without}} = \text{softmax}(L_T^{\text{without}})$ and $P^{\text{with}} = \text{softmax}(L_T^{\text{with}})$, using three metrics:

- 1242 • **KL Divergence (D_{KL}):** Measures how much the output probability distribution of the target
 1243 module changes when the source module's information is added. A larger value indicates a greater
 1244 change.

$$D_{KL} = D_{KL}(P^{\text{with}} \parallel P^{\text{without}}) \quad (27)$$

- 1245 • **Prediction Change Rate (Δ_{pred}):** The fraction of samples in a batch for which the predicted class
 1246 (the one with the highest probability) changes. $\mathbb{I}(\cdot)$ is the indicator function.

$$\Delta_{\text{pred}} = \mathbb{E} [\mathbb{I}(\arg \max(L_T^{\text{with}}) \neq \arg \max(L_T^{\text{without}}))] \quad (28)$$

- 1247 • **Confidence Change (Δ_{conf}):** The absolute difference in the average prediction confidence (the
 1248 value of the max probability) between the two conditions.

$$\Delta_{\text{conf}} = |\mathbb{E}[\max(P^{\text{with}})] - \mathbb{E}[\max(P^{\text{without}})]| \quad (29)$$

1249 **Final Influence Score** The final Information Flow Influence Score (IFS) for the connection $M_S \rightarrow$
 1250 M_T is the simple average of these three normalized metrics, providing a single value to represent the
 1251 strength of the connection.

$$\text{IFS}(M_S \rightarrow M_T) = \frac{1}{3}(D_{KL} + \Delta_{\text{pred}} + \Delta_{\text{conf}}) \quad (30)$$

A.9 SPEECH UNDERSTANDING AND REASONING EVALUATION

The prompt used to have GPT-4o score the reasoning and responses generated by our models is as follows:

1266 **System Prompt:**

1267 You are a strict model-as-judge. You will assign two scores on a 0-10 scale for each sample: one
 1268 for the model's reasoning ("analysis") and one for the model's response ("response").
 1269 You are given: the ASR transcript, the final.STATE with four labels (WMA, ToM, SA, Prag), plus the
 1270 sample's analysis and response.

1271 You must follow the rubric and output strict JSON only. No extra explanations, no code blocks, no
 1272 markdown – only valid JSON.
 1273 [Scoring Rubric]

1274 (1) reasoning_score (0-10) – for analysis:

- Completeness: does it cover and explain all four dimensions (WMA, ToM, SA, Prag) and their relations? Does it reference key ASR info?
- Accuracy: does it align with final.STATE labels exactly? No inventing or changing labels.
- Causal logic: does it present a clear, coherent causal chain from ASR → state judgment → pragmatic intent?

1275 Guide:

1276 10: all four dimensions, fully accurate, causal and clear.

1277 8: mostly correct, minor gaps or generalization.

1278 6: partially correct, superficial, 1-2 inaccuracies.

1279 4: several inaccuracies or weak logic, clear mismatches.

1280 2: mostly irrelevant or confused reasoning.

1281 0: missing, fabricated, or totally wrong.

1282 (2) response_score (0-10) – for response:

- Logicality and appropriateness: does it respond reasonably to the ASR context (e.g. anger, wh-question)? Polite, safe wording?
- Match with final.STATE: does it reflect WMA, ToM, SA, Prag consistently? No contradictions.

1283 Guide:

1284 10: highly consistent, natural, well-suited strategy.

1285 8: generally consistent, minor flaws or templated.

1286 6: partially consistent, weak on key info.

1287 4: multiple mismatches, awkward.

1288 2: mostly irrelevant or inappropriate.

1289 0: nonsensical, unsafe, or severely inconsistent.

1290 EXAMPLE:

1291 [Output Format – strict JSON only]

1292 {

1293 "reasoning_score": <float 0-10>,

1294 "response_score": <float 0-10>

1295 }

1296 A.10 THE USE OF LARGE LANGUAGE MODELS
12971298 During the writing of this paper, we leveraged GPT-5 to aid in polishing the writing. Its use was
1299 exclusively limited to improving grammar, phrasing, and overall readability. The core scientific
1300 contributions, experimental results, and intellectual content are entirely our own.
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349