

# STATE-SPACE MODELING IN NATURAL LANGUAGE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In many real-world applications, the goal is to infer the latent dynamics of an underlying process from a sequence of observations. Traditionally, state-space models such as hidden Markov models (HMMs) have been used to represent dynamic systems by positing a discrete or continuous latent state space, where states evolve according to a *transition* model and observations follow an *emission* model. However, these models are typically limited to structured observations and numerical state representations that are often not interpretable or identifiable. In this paper, we propose a new class of state-space models that formulates *state inference* and *prediction* as language tasks, with both the transition and emission models implemented by a large language model (LLM). In our model, latent states are represented as concise natural-language descriptions, while observations may comprise large and unstructured text corpora. We develop a post-training procedure to fine-tune LLMs for state inference and prediction, inspired by variational inference algorithms in classical state-space modeling. We demonstrate the promise of this approach through preliminary experiments on disease progression modeling from clinical notes and on modeling geopolitical relationships using news articles.

**Track:** Research

## 1 INTRODUCTION

Large language models (LLMs) have limited context windows, and their performance can degrade sharply as prompt length increases (Hong et al., 2025). In many real-world applications, the input context is produced by an evolving process that continually accumulates text over time. For example, in healthcare settings, clinicians may wish to use LLMs to assist in analyzing a patient’s expanding electronic health record (EHR) (Sellergren et al., 2025). Fundamentally, patients in an EHR dataset can be viewed as independent realizations of stochastic health trajectories governed by shared latent dynamics. In such settings, a potentially effective strategy for compressing the growing input context is to model these latent dynamics, thereby learning to extract past information that is most relevant at each time point and most useful for downstream predictions at future time points.

A variety of methods have been proposed for handling long-context inputs, yet these approaches generally do not address settings in which the context itself forms a time series. Previous work includes methods based on sparse attention (Roy et al., 2021), memory-compression architectures such as Compressive Transformers (Rae et al., 2019), and inference-time methods including context condensation and recursive language modeling (Khattab et al., 2021; Zhang et al., 2025). In this paper, however, we do not focus on long-context modeling in general, but rather on contexts generated by real-world, time-stamped stochastic processes that produce continually growing text corpora. Examples include sequences of clinical notes in a patient’s EHR, which document patient visits over time, as well as evolving collections of news articles describing geopolitical events or financial markets. Our key hypothesis is that learning to characterize the shared latent dynamics underlying these processes is equivalent to learning how to compress their growing context at each point in time.

Building on the observations above, we propose a state-space modeling approach for processing long contexts generated by temporal processes. State-space models (SSMs) have long served as a foundational framework for modeling sequential data and dynamical systems in control theory (Hamilton, 1994), signal processing (Rao & Arun, 2002), and time-series analysis (Aoki, 2013; Durbin & Koopman, 2012). In an SSM, the evolution of an observed sequence is governed by an underlying latent state with structured *transition* dynamics, while observations are generated through a separate *emission* process. This separation between latent dynamics and observations yields a compact and

054 interpretable representation of temporal dependencies, which enables SSMs to capture long-range  
 055 structure through recursive state updates rather than explicit pairwise interactions. Previous work has  
 056 extended SSMs using deep neural networks to parameterize state transitions and emission distribu-  
 057 tions, as in the case of Deep Markov Models and its variants (Krishnan et al., 2015; 2017; Johnson  
 058 et al., 2016; Alaa & van der Schaar, 2019). However, these approaches have largely been limited to  
 059 structured or tabular data in traditional time-series modeling settings.

060 Our modeling approach conceptualizes context compression as a state-space modeling problem. In  
 061 our framework, sequences of text corpora are assumed to be generated by an emission distribution  
 062 conditioned on an underlying latent state representing the context, while state evolution satisfies a  
 063 Markov property in which each state depends only on its predecessor. Crucially, these latent states  
 064 are represented as natural-language summaries of the observed text at each time step. Prior work sug-  
 065 gests that training objectives based on predicting future summaries can outperform standard next-  
 066 token prediction in LLMs (Mahajan et al., 2025). From this perspective, context compression be-  
 067 comes equivalent to modeling the dynamics of an evolving system rather than truncating or heuris-  
 068 tically condensing input. This state-space decomposition allows us to formulate both *state inference*  
 069 (context compression) and *state prediction* as language modeling tasks.

070 To operationalize our proposed SSM-based framework, we develop a heuristic post-training proce-  
 071 dure for fine-tuning LLMs to perform state inference and state prediction, drawing inspiration from  
 072 variational inference methods in classical state-space modeling. We demonstrate the effectiveness  
 073 and utility of this approach through preliminary experiments on disease progression modeling using  
 074 clinical notes and on modeling geopolitical relationships from news articles.

## 076 2 METHODS

078 We consider the problem of sequential modeling over a time-indexed sequence of text corpora

$$079 X_{1:T} = (X_1, X_2, \dots, X_T), \quad X_t \in \mathcal{X},$$

081 where each  $X_t$  denotes observed, unstructured text (e.g., news articles or clinical notes) at time step  
 082  $t$ . The collection  $X_{1:T}$  constitutes the available context at time  $T$ . A common objective is to perform  
 083 a downstream task by responding to a query  $Q$  conditioned on this context using a language model,  
 084 i.e.,  $Y = \text{LLM}(Q; X_{1:T})$ . For instance, the query  $Q$  may be a question about a patient’s likely health  
 085 outcomes in the future given all the clinical notes  $X_{1:T}$  collected in their previous visits. However,  
 086 performance on such tasks often degrades as the context  $X_{1:T}$  grows, which motivates the need for  
 087 principled methods to capitalize on the temporal structure of  $X_{1:T}$  to compress the context.

088 **State-space modeling of temporal context.** We assume that the sequence  $X_{1:T}$  is generated via a  
 089 state-space model (SSM), with a factorized joint distribution of observations and latent states:

$$091 p(X_{1:T}, Z_{1:T}) = p(Z_1) \prod_{t=2}^T p(Z_t | Z_{t-1}) \prod_{t=1}^T p(X_t | Z_t), \quad (1)$$

094 where  $Z_t \in \mathcal{Z}$  is the latent state of the dynamic process generating document  $X_t$  at time  $t$ . Here, we  
 095 assume that this process is Markovian: the states evolve so that each new state  $Z_t$  depends only on  
 096 the previous state  $Z_{t-1}$ . Thus, the joint distribution of observed text and latent states can be fully de-  
 097 termined by an *emission* distribution  $p(X_t | Z_t)$  and a *state transition* distribution  $p(Z_t | Z_{t-1})$ .

098 We further assume that the state space  $\mathcal{Z}$  consists of natural-language strings that explicitly encode  
 099 the system’s underlying state. This enables latent states to be both expressive and interpretable, while  
 100 maintaining compatibility with standard language modeling objectives. For example, in an EHR set-  
 101 ting, a latent state may capture a patient’s evolving health across multiple comorbidities summarized  
 102 in textual form. Individual observations  $X_t$  (e.g., physician notes) are then generated as condition-  
 103 ally independent views of this state. This perspective naturally accommodates long contexts, as in-  
 104 formation is propagated through the recurrent state dynamics rather than requiring direct interactions  
 105 among all past observations. Moreover, representing states in natural language allows the model to  
 106 leverage pretrained language priors to derive principled algorithms for state inference.

107 **Learning and inference.** By representing latent states as natural-language strings, we can formulate  
 both learning and inference as standard language modeling tasks. Concretely, the transition model

108  $p(Z_t|Z_{t-1})$  and emission model  $p(X_t|Z_t)$  can be parameterized as autoregressive language models,  
 109 where state evolution and observation generation can be learned via next-token prediction objectives.  
 110 Inference over latent states similarly reduces to conditional text generation. Given prior context, the  
 111 posterior distribution over the latent states  $p(Z_t|X_{1:t})$  can be approximated by prompting a language  
 112 model to generate a state description consistent with past states and current observations.

113 A standard way to conduct learning in inference with deep learning models is to jointly train (i) an inference  
 114 model and (ii) the emission and transition models by maximizing the marginal log-likelihood  
 115  $\log p(X_{1:T})$  end-to-end (Krishnan et al., 2017). This objective is generally intractable, therefore it  
 116 is common to optimize its variational lower bound (ELBO) with an amortized inference model:

$$117 \mathcal{L}_{\theta, \phi} = \sum_{t=1}^T \log p_{\text{LLM}_{\theta}}(X_t|Z_t) + \log p_{\text{LLM}_{\theta}}(Z_1) + \sum_{t=2}^T \log p_{\text{LLM}_{\theta}}(Z_t|Z_{t-1}) - \log q_{\text{LLM}_{\phi}}(Z_{1:T}|X_{1:T}),$$

120 where  $p_{\text{LLM}_{\theta}}$  and  $q_{\text{LLM}_{\phi}}$  are the token distributions induced by two LLMs with parameters  $\theta$  and  $\phi$ ,  
 121 or equivalently, a single LLM steered to model observations, states as well as state inferences using  
 122 different prompts. We developed a simplification of the ELBO loss that first starts by labeling the  
 123 latent states  $\{Z_t\}$  using a high-capacity teacher model applied with backward-induction, and then  
 124 training the inference model in a supervised fashion. Full pseudocode is provided in the Appendix  
 125 (Algorithm 1).

## 126 2.1 BASELINES AND EVALUATION

127 We compare against embedding-based HMMs, TopicGPT+HMM, Neural HMMs, a forward LLM  
 128 baseline, and an LSTM classifier (see Appendix A.2). Discrete baselines use latent cluster indices as  
 129 states. Textual summaries from the Forward LLM and NL-SSM are mapped to binary regimes using  
 130 Goldstein scores (Goldstein, 1992; Gerner et al., 2009). We report Normalized Mutual Information  
 131 (NMI) for geopolitical regime alignment, and sensitivity/F1 for CKD progression detection.

## 132 3 RESULTS

### 133 3.1 GEOPOLITICAL REGIME DISCOVERY

134 We evaluate on the Global Database of Events, Language, and Tone (GDELT) dataset (Leetaru &  
 135 Schrodt, 2013), a global database of timestamped news events. Fixed sentence embeddings (B1) and  
 136 neural HMMs (B3) failed to discover meaningful regime structure ( $\text{NMI} \leq 0.01$ ; see Table 1), cluster-  
 137 ing by geography rather than regime. TopicGPT achieved greater model alignment ( $\text{NMI}: 0.197$ ),  
 138 demonstrating that LLM-induced topic abstractions provide effective state descriptions. However,  
 139 overlapping conflict and cooperation signals across topics suggest static assignments conflate mixed  
 140 regimes.

141 Table 1: Normalized Mutual Information (NMI) between inferred states and ground-truth regimes.

Method	NMI $\uparrow$
B1: Fixed Embed + HMM	0.006
B2: TopicGPT + HMM	0.197
B3: Neural HMM	0.002
B4: Forward LLM	0.129
<b>NL-SSM</b>	<b>0.551</b>

142 **NL-SSM discovers meaningful geopolitical regimes.** NL-SSM achieved substantially higher  
 143 alignment ( $\text{NMI}: 0.551$ ), producing temporally discriminative state summaries. In contrast, the  
 144 Forward LLM relied on hedging language (e.g., would, possibly, could, might), and frequently predicted  
 145 persistent cooperation.

### 146 TEMPORAL DYNAMICS OF LEARNED STATES

147 We quantified temporal variation by computing cosine similarity between consecutive states using  
 148 sentence-transformer embeddings. As shown in Fig. 1, the Forward LLM shows strong persistence  
 149  
 150  
 151  
 152  
 153  
 154  
 155  
 156  
 157  
 158

on GDELT (95% above 0.5; mean=0.71, std=0.12), whereas NL-SSM states vary more (42% above 0.5; mean=0.46, std=0.32;  $p < 0.0001$ ). CKD shows the same trend: Forward LLM summaries remain nearly identical (100%; mean=0.93, std=0.06), while NL-SSM evolves across patient visits (71%; mean=0.68, std=0.36;  $p < 0.0001$ ), suggesting future-conditioned supervision encourages temporally adaptive states.

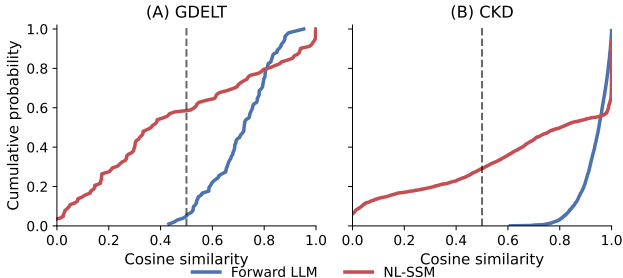


Figure 1: Cumulative distribution of cosine similarity between consecutive states on GDELT and CKD datasets.

### 3.2 CLINICAL STATE INFERENCE (NEPHROLOGY)

#### CLINICAL CONTENT ANALYSIS

We evaluate NL-SSM on chronic kidney disease (CKD) prognosis (13.2% progression rate). We analyze clinical content with nephrology-informed categories; see Appendix A.8 and Table 3. NL-SSM summaries contain significantly more objective clinical evidence (lab values, stage mentions, treatments, comorbidities;  $p < 0.0001$ ), whereas the Forward LLM relies more on speculative trajectory and prognosis language. This mirrors the persistence bias observed in GDELT, and suggests future-conditioned training yields physiologically grounded state summaries.

#### CKD PROGRESSION

NL-SSM achieves higher sensitivity (0.63 vs 0.38) and F1 (0.25 vs 0.21) than Forward LLM (Table 2), but both exhibit lower precision (15%), with false alarm rates of 52% (NL-SSM) and 34% (Forward LLM). This precision-sensitivity trade-off reflects the challenge of rare-event detection: improving recall for critical transitions comes at the cost of increased false alarms. A qualitative CKD progression example is shown in Appendix 2.

Table 2: CKD progression detection performance on all test samples ( $N = 2567$ ). Missing predictions count as incorrect. 95% confidence intervals via patient-clustered bootstrap.

Method	F1	Sensitivity	Coverage
Forward LLM	0.21 [0.18, 0.25]	0.38 [0.33, 0.43]	81.6%
NL-SSM	<b>0.25 [0.22, 0.28]</b>	<b>0.63 [0.56, 0.69]</b>	<b>96.0%</b>

## 4 DISCUSSION

Our results demonstrate that training with future context can reshape the learned state space. While embedding-based and topic-based probabilistic state-space model baselines cluster by static semantic features and the Forward LLM exhibits persistence bias, NL-SSM produces time-specific state descriptions that better capture temporal evolution. This highlights a fundamental tradeoff in persistent regimes: stability optimization achieves high accuracy but systematically misses critical transitions.

**Limitations.** Evaluating natural language state summaries is challenging. While projecting to categorical labels enables quantitative evaluation, it may not capture nuanced information. Second, NL-SSM requires a high-capacity teacher during training, increasing computational cost.

## 216 REFERENCES

- 217
- 218 Ahmed M Alaa and Mihaela van der Schaar. Attentive state-space modeling of disease progression.  
219 *Advances in neural information processing systems*, 32, 2019.
- 220 Masanao Aoki. *State space modeling of time series*. Springer Science & Business Media, 2013.
- 221
- 222 James Durbin and Siem Jan Koopman. *Time series analysis by state space methods*. Oxford univer-  
223 sity press, 2012.
- 224 Deborah J. Gerner, Philip A. Schrodtt, Omur Yilmaz, and Erin Weddle. *Conflict and Mediation*  
225 *Event Observations (CAMEO) Manual*, 2009. URL [http://data.gdeltproject.org/](http://data.gdeltproject.org/documentation/CAMEO.Manual.1.1b3.pdf)  
226 [documentation/CAMEO.Manual.1.1b3.pdf](http://data.gdeltproject.org/documentation/CAMEO.Manual.1.1b3.pdf). Version 1.1b3.
- 227
- 228 Joshua S Goldstein. A conflict-cooperation scale for weis events data. *Journal of conflict resolution*,  
229 36(2):369–385, 1992.
- 230 James D Hamilton. State-space models. *Handbook of econometrics*, 4:3039–3080, 1994.
- 231
- 232 Kelly Hong, Anton Troynikov, and Jeff Huber. Context rot: How increasing input tokens impacts  
233 llm performance. URL <https://research.trychroma.com/context-rot>, retrieved October, 20:2025,  
234 2025.
- 235 Matthew Johnson, David K Duvenaud, Alexander Wiltschko, Ryan P Adams, and Sandeep R Datta.  
236 Composing graphical models with neural networks for structured representations and fast infer-  
237 ence. In *NeurIPS*, 2016.
- 238
- 239 Omar Khattab, Christopher Potts, and Matei Zaharia. Baleen: Robust multi-hop reasoning at scale  
240 via condensed retrieval. *Advances in Neural Information Processing Systems*, 34:27670–27682,  
241 2021.
- 242 Rahul G Krishnan, Uri Shalit, and David Sontag. Deep kalman filters. *arXiv preprint*  
243 *arXiv:1511.05121*, 2015.
- 244
- 245 Rahul G Krishnan, Uri Shalit, and David Sontag. Structured inference networks for nonlinear state  
246 space models. In *AAAI*, 2017.
- 247 Kalev Leetaru and Philip A. Schrodtt. Gdelt event database codebook v2.0, 2013.  
248 URL [https://data.gdeltproject.org/documentation/GDELT-Event\\_](https://data.gdeltproject.org/documentation/GDELT-Event_Codebook-V2.0.pdf)  
249 [Codebook-V2.0.pdf](https://data.gdeltproject.org/documentation/GDELT-Event_Codebook-V2.0.pdf).
- 250
- 251 Divyat Mahajan, Sachin Goyal, Badr Youbi Idrissi, Mohammad Pezeshki, Ioannis Mitliagkas, David  
252 Lopez-Paz, and Kartik Ahuja. Beyond multi-token prediction: Pretraining llms with future sum-  
253 maries. *arXiv preprint arXiv:2510.14751*, 2025.
- 254 Chau Minh Pham, Alexander Hoyle, Simeng Sun, Philip Resnik, and Mohit Iyyer. Topicgpt: A  
255 prompt-based topic modeling framework. doi: 10.48550. *arXiv preprint arXiv:2311.01449*, 2024.
- 256
- 257 Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, and Timothy P Lillicrap. Compressive  
258 transformers for long-range sequence modelling. *arXiv preprint arXiv:1911.05507*, 2019.
- 259 Bhaskar D Rao and KS Arun. Model based processing of signals: A state space approach. *Proceed-*  
260 *ings of the IEEE*, 80(2):283–309, 2002.
- 261
- 262 Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse  
263 attention with routing transformers. *Transactions of the Association for Computational Linguis-*  
264 *tics*, 9:53–68, 2021.
- 265 Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo  
266 Kohlberger, Shawn Xu, Fayaz Jamil, Cian Hughes, Charles Lau, et al. Medgemma technical  
267 report. *arXiv preprint arXiv:2507.05201*, 2025.
- 268
- 269 Alex L Zhang, Tim Kraska, and Omar Khattab. Recursive language models. *arXiv preprint*  
*arXiv:2512.24601*, 2025.

270 A APPENDIX

271  
272 A.1 FULL NL-SSM ALGORITHM

---

273  
274 **Algorithm 1** Natural Language State-Space Model (NL-SSM)

---

275 **Require:** Document sequence  $\{x_t\}_{t=1}^T$   
 276 **Ensure:** Latent state sequence  $\{s_t\}_{t=1}^T$

277  
278 1: **Backward label induction (teacher, offline)**

279 2:  $s_T \leftarrow \text{LM}_{\text{teach}}(x_T)$

280 3: **for**  $t = T - 1$  **to** 1 **do**

281 4:      $s_t \leftarrow \text{LM}_{\text{teach}}(x_t, s_{t+1})$

282 5: **end for**

283  
284 6: **Train state summarization model**

285 7: **for**  $t = 1$  **to**  $T$  **do**

286 8:     Train  $\text{LM}_{\text{summ}}$  on  $(s_{t-1}, x_t) \rightarrow s_t$

287 9: **end for**

288  
289 10: **Train state transition model**

290 11: **for**  $t = 1$  **to**  $T - 1$  **do**

291 12:     Train  $\text{LM}_{\text{trans}}$  on  $s_t \rightarrow s_{t+1}$

292 13: **end for**

293 14: **Inference (causal filtering)**

294 15:  $\hat{s}_1 \leftarrow \text{LM}_{\text{summ}}(x_1)$

295 16: **for**  $t = 2$  **to**  $T$  **do**

296 17:      $\hat{s}_t \leftarrow \text{LM}_{\text{summ}}(\hat{s}_{t-1}, x_t)$

297 18: **end for**

298  
299 19: **Forecasting**

300 20: **for**  $t = 1$  **to**  $T - 1$  **do**

301 21:      $\tilde{s}_{t+1} \leftarrow \text{LM}_{\text{trans}}(\hat{s}_t)$

302 22: **end for**

---

303  
304  
305 A.2 BASELINES

306 **Baseline 1: Fixed Embeddings + HMM.** We embed each document using a frozen  
 307 all-MiniLM-L6-v2 encoder, and form monthly observations averaging these embeddings. A  
 308 Gaussian-emission HMM is then trained using Expectation-Maximization to maximize marginal  
 309 likelihood. HMMs are trained with the expectation-maximization algorithm until convergence of  
 310 the log-likelihood.

311 **Baseline 2: TopicGPT + HMM.** We adopt TopicGPT Pham et al. (2024), which uses LLMs to  
 312 induce topics from text corpora. We use GPT-4o to produce a global topic list (see Appendix A.3)  
 313 from 100 sampled training articles, assign topics to three randomly sampled articles per month per  
 314 country pair, and train a categorical HMM. Topics are fixed across country-pair splits.

315 **Baseline 3: Neural HMM.** This baseline jointly trains a neural encoder and an HMM using Gener-  
 316 alized EM. The encoder  $f_\phi$  maps each document to a continuous embedding and is instantiated  $f_\phi$   
 317 using all-MiniLM-L6-v2 with only the final transformer layer fine-tuned.

318 **Baseline 4: Forward LLM.** We use the same high-capacity language model (GPT-4o) as a purely  
 319 causal baseline. At each timestep, the model produces a state summary conditioned only on the  
 320 previous state summary (if available) and the current observations (three randomly sampled articles).  
 321 The model is prompted with the same task description as NL-SSM (see Appendix A.4), but is never  
 322 trained or conditioned on future states. This baseline isolates the effect of future-conditioned training  
 323 supervision while controlling for model capacity.

**Baseline 5: Long Short-Term Memory (LSTM).** We train an LSTM classifier to predict regime labels from `all-MiniLM-L6-v2` embeddings using a 6-month sliding window. Parameters are shared across country pairs.

**Our Contribution: NL-SSM.** We use GPT-4o as a teacher model during training to generate future-conditioned state summaries, and Phi-3-mini-128k as student models for state summarization, and state transition. During inference, NL-SSM operates strictly causally, and observes the same information as the Forward LLM baseline (i.e., only the previous timestep’s inferred state and current observation), isolating the effect of future-conditioned training supervision (see Appendix A.5).

### A.3 BASELINE 2: TOPICGPT

TopicGPT injects task-specific guidance by explicitly constraining the language model, GPT 4o, to produce relation-level abstractions (e.g., diplomatic, military, or economic interactions) rather than generic semantic topics. However, this guidance is static and corpus-level: topics are induced once from the training data, do not condition on temporal context, and do not incorporate future information. As a result, TopicGPT captures coarse relational semantics, but is not designed to model temporal regime transitions. We produce a fixed set of reusable, relation-level topics that characterize bilateral country relations across all country pairs and time periods.

Topics are induced using only documents from the training split to avoid information leakage. We first restrict the corpus to training examples and randomly sample 100 documents for topic induction using a fixed random seed to ensure reproducibility.

**Prompt:** The language model is prompted as follows:

You are analyzing how relations between two countries evolve over time.

**Context:**

You are given a set of real-world news articles.

**Task:**

Produce 20 topics describing the bilateral relations between the countries during this time period.

**Rules:**

- Each topic must be a short label (no description)
- Labels must be 1–3 words
- Topics must be reusable across all country pairs

**Return format:**

```
topic_list:
- topic1
- topic2
...
```

The final global topic list consists of the following 20 relation-level categories:

- Military Cooperation
- Diplomatic Tensions
- Trade Relations
- Humanitarian Aid
- Political Protests
- Security Concerns
- Economic Sanctions
- Cultural Exchange
- Leadership Changes
- Peace Negotiations

- 378 • Territorial Disputes
- 379 • Election Influence
- 380 • Energy Partnerships
- 381 • Technological Collaboration
- 382 • Immigration Policies
- 383 • Defense Alliances
- 384 • Crisis Management
- 385 • Environmental Agreements
- 386 • Human Rights Issues
- 387 • Strategic Partnerships

#### 391 A.4 BASELINE 4: FORWARD REASONING LLM PROMPTS (GDELT)

392 The Forward Reasoning LLM baseline is a strictly causal approach that mirrors the inference-time  
 393 structure of NL-SSM, but without future-conditioned supervision. The model operates in two stages  
 394 at each timestep: (1) inferring the current latent state using observed text, and (2) forecasting the  
 395 next latent state without access to future observations.  
 396  
 397

398 **Prompt 1: State Inference with Observations.** At each timestep  $t$ , the language model infers  
 399 the current state summary conditioned on the previous state summary (if available) and the current  
 400 month’s news articles.  
 401

402 You are analyzing how relations between two countries evolve over time.

403 **Task:**

404 Based only on the prior month’s summary (if available) and the current articles, write a  
 405 concise 2–3 sentence summary describing the current bilateral relations between the coun-  
 406 tries.

407 **Country pair:**  $[A-B]$

408 **Time period:**  $[Month, Year]$

409 **Prior state summary (from the previous month):**

410  $[Previous\ state\ summary]$

411 If no prior state is available, the model is explicitly informed that no prior summary exists.

412 **Current articles:**

413 *Article 1: [Text]*

414 *Article 2: [Text]*

415 *Article 3: [Text]*

416 **State summary:**

417  
 418 **Prompt 2: State Transition.** To evaluate temporal prediction, the model is then prompted to  
 419 forecast the next state summary using only the inferred state at time  $t$ , without access to any articles  
 420 from time  $t+1$ .  
 421

422 You are analyzing how relations between two countries evolve over time.

423 **Task:**

424 Based only on the current state summary, predict the bilateral relations at the next time  
 425 period. Write a concise 2–3 sentence summary describing the expected state and trajectory.

426 **Country pair:**  $[A-B]$

427 **Current time period:**  $[Month\ t]$

428 **Next time period:**  $[Month\ t+1]$

429 **Current state summary:**

430  $[State\ summary\ at\ time\ t]$

431 **Predicted next state summary:**

**Decoding Details.** All forward baseline prompts use greedy decoding with temperature set to 0 and a maximum response length of 1500 tokens. The model is queried strictly causally and never conditions on future observations.

#### A.5 NL-SSM TRAINING AND INFERENCE PROMPTS

**Train 1 (State Inference).** Trains  $(\text{teacher\_state}_{t-1}, x_t) \rightarrow \text{teacher\_state}_t$ .

**Task:** Using the prior month’s summary (if available) and current articles, write a 2–3 sentence summary of current bilateral relations.

**Country pair:** [A–B], Time: [Month  $t$ ]

There is no prior state summary *or*

Prior state summary (from the previous month): [teacher\\_state $_{t-1}$ ]

Articles: Article 1, Article 2, Article 3

**State summary:**

**Train 2 (State Transition).** Trains  $\text{teacher\_state}_t \rightarrow \text{teacher\_state}_{t+1}$ .

**Task:** Using only the current state summary, predict relations in the next time period.

**Country pair:** [A–B], Time: [Month  $t$ ]

**Current state summary:** [teacher\\_state $_t$ ]

**Next state summary:**

**Inference.** Two-pass decoding: (i) infer states using Train 1 prompt and observed articles; (ii) forecast next states using Train 2 prompt without observations.

#### A.6 IMPLEMENTATION DETAILS.

**Teacher Model (Backward Label Induction):** GPT-4o via OpenAI API. Temperature=0, max tokens=1500.

**Student Models:** Phi-3-mini-128k fine-tuned with LoRA ( $r=16$ ,  $\alpha=32$ , dropout=0.05). Training: 3 epochs, learning rate=2e-4, effective batch size=16 (per-device=1, gradient accumulation=16).

**Evaluation:** Sentence-BERT (all-MiniLM-L6-v2) for computing state similarity. Goldstein score mapping for regime classification.

#### A.7 CKD CLINICAL MARKER LEXICON

We analyze clinical content in predicted CKD states using the following marker categories. Terms with asterisks (\*) use prefix matching (e.g., “declin” matches “decline,” “declining,” “declined”).

- **Lab values:** egfr, creatinine, gfr, bun, potassium, phosphorus, hemoglobin
- **Stage mentions:** stage 1, stage 2, stage 3, stage 4, stage 5, esrd
- **Trajectory:** progress, declin\*, stable, improv\*, worsen\*, slow, rapid
- **Treatment:** dialysis, transplant, medication, ace, arb, diuretic, epo
- **Comorbidities:** diabetes, hypertension, anemia, proteinuria, edema
- **Prognosis:** risk, prognosis, outlook, future, expect\*, likely

*Note:* All matching is case-insensitive. Medical abbreviations: ace = ACE inhibitors, arb = angiotensin receptor blockers, epo = erythropoietin.

## A.8 CLINICAL CONTENT ANALYSIS

Table 3: Clinical content in CKD state predictions (mean mentions per summary). \*\*\* $p < 0.001$ .

Category	NL-SSM	Forward LLM
Lab values***	<b>2.08</b>	1.81
Stage mention***	<b>1.69</b>	0.95
Comorbidities***	<b>0.91</b>	0.76
Treatment ***	0.77	0.46
Trajectory***	1.81	<b>2.31</b>
Prognosis***	1.53	<b>2.42</b>

## A.9 CKD QUALITATIVE EXAMPLE

## Qualitative Example: CKD Progression (Stage 3b → Stage 4)

## NL-SSM

## Forward LLM

State at time  $t$  (Input)

Stage 3 chronic kidney disease (CKD) with an estimated glomerular filtration rate (eGFR) of 37 mL/min/1.73 m<sup>2</sup>...**expected to progress to stage 4 CKD, with a high risk of end-stage renal disease (ESRD) in the future, necessitating dialysis or kidney transplantation.**

Stage 3 chronic kidney disease (CKD), with an eGFR of 37 ml/min/1.73 m<sup>2</sup> and a creatinine level of 1.3...**indicates a slight decline in kidney function**, with worsening edema and shortness of breath potentially related to fluid overload, necessitating the initiation of diuretic therapy.

Forecasted state at time  $t+1$  (Output)

**Stage 4 CKD** with an estimated glomerular filtration rate (eGFR) of 25 mL/min/1.73 m<sup>2</sup>, **primarily due to hypertensive nephrosclerosis**. The disease trajectory is progressive, with a high risk of advancing to end-stage renal disease (ESRD) in the future, **necessitating dialysis or kidney transplantation.**

**Stage 3 CKD**, with a potential slight decline in eGFR to around 35 ml/min/1.73 m<sup>2</sup>, given the current trajectory of kidney function. The **prognosis will focus on continued management of hypertension and fluid overload**, with adjustments to diuretic therapy as needed.

Figure 2: Qualitative comparison of forecasts before a true CKD stage transition. Both models receive comparable patient summaries at time  $t$  (gray). NL-SSM correctly anticipates progression to Stage 4 at  $t+1$  (green), whereas the Forward LLM predicts persistence at Stage 3 (red).