

From Word Sequences to Behavioral Sequences: Adapting Modeling and Evaluation Paradigms for Longitudinal NLP

Anonymous ACL submission

Abstract

While NLP typically treats documents as independent and unordered samples, in longitudinal studies, this assumption rarely holds: documents are nested within authors and ordered in time, forming person-indexed, time-ordered *behavioral sequences*. Here, we demonstrate the need for and propose a longitudinal modeling and evaluation paradigm that consequently updates four parts of the NLP pipeline: (1) evaluation splits aligned to generalization over people (*cross-sectional*) and/or time (*prospective*); (2) accuracy metrics separating between-person differences from within-person dynamics; (3) sequence inputs to incorporate history by default; and (4) model internals that support different *coarseness* of latent state over histories (pooled summaries, explicit dynamics, or interaction-based models). We demonstrate the issues ensued by traditional pipeline and our proposed improvements on a dataset of 17k daily diary transcripts paired with PTSD symptom severity from 238 participants, finding that traditional document-level evaluation can yield substantially different and sometimes reversed conclusions compared to our ecologically valid modeling and evaluation. We tie our results to a broader discussion motivating a shift from word-sequence evaluation toward *behavior-sequence* paradigms for NLP.

1 Introduction

NLP typically frames prediction as mapping isolated instances of language to an outcome. However, recent work increasingly targets inherently longitudinal problems, where multiple instances of writing are accrued per person (Kumar and Carley, 2019; Sawhney et al., 2021; V Ganesan et al., 2021; Tsakalidis et al., 2022). In these regimes, documents are generated by *people* over *time*, yielding person-indexed, time-ordered behavioral sequences rather than independent instances that current standard statistical techniques are designed for. This

mismatch propagates through the standard NLP pipeline: random document splits can leak person-specific signal and scramble temporal order, per-document training ignores informative history, and flattened accuracy metrics lack clear real-world generalization goals. As a result, models can appear stronger than they really are for realistic applications and conclusions about the relationship between language and outcomes can even be flipped (e.g. see Figure 1).

These violations of the independence assumption are not rare; they arise when many documents are produced or labeled by a limited set of authors (Geva et al., 2019), and they are further amplified in longitudinal data where documents are also ordered in time. In such “human-level” applications – such as predicting mental health, personality, or demographics from text – the effective sample size is not the number of *documents* in a dataset but the number of *individuals* (V Ganesan et al., 2021), and variation decomposes into between-person differences and within-person change (Curran and Bauer, 2011; Hoffman and Stawski, 2009). However traditionally used flattened, document-level metrics conflate these components, rewarding models that memorize stable person idiosyncrasies and obscuring whether they capture within-person dynamics (Hamaker et al., 2015).

These observations motivate an ontological shift in NLP pipeline: treat documents as ordered samples of *verbal behavior* emitted by individuals, not as word sequences (Boyd and Schwartz, 2021; V Ganesan et al., 2024; Soni et al., 2024). Formally, each observation is indexed by person i and time t , $(x_{i,t}, y_{i,t})$, with dependence induced by (a) person-level structures (style, baseline levels) and (b) temporal structures (autocorrelation, trends). Consequently, evaluation metrics should distinguish whether performance is driven by generalization to stable person-level differences, or dynamic within-person changes.

When documents are nested within individuals and ordered in time, evaluation must specify what is held out (people, time, or both), since each choice tests a different generalization problem. We therefore separate two axes: (1) *cross-sectional generalization*: does a model generalize to unseen people? and (2) *prospective generalization*: does a model generalize forward in time for the same people? These targets are not interchangeable, and collapsing them into a single document-level evaluation can yield misleading inferences about model usefulness.

Here, we introduce a modeling and evaluation framework for longitudinal ML/NLP. We show why such a framework is necessary when text is treated as a repeated measurement of a person over time. Rather than optimizing a single leaderboard endpoint, we ask what a model’s performance can legitimately claim under different deployment-relevant generalization targets. We demonstrate these points on a rich longitudinal dataset, showing how standard NLP/ML evaluation choices can mischaracterize performance when documents are treated as independent examples.

We make five contributions, organized from problem to implications: (1) Demonstrating the problem: We show that traditional document-level evaluation can yield substantially different conclusions than evaluations tied to ecologically valid, real-world use-cases. (2) Developing an evaluation framework: We operationalize evaluation splits that target generalization across *people* (cross-sectional) and across *time* (prospective), and show that these settings support qualitatively different inferences. (3) Demonstrate need for metric decomposition: We propose reporting both *between-person* and *within-person* variants of standard metrics to disentangle person-level and temporal performance. (4) Show the need to adapt modeling: Incorporating temporal context improves performance under these changes and brings out longitudinal modeling challenges. (5) Broader implications: We connect these findings to a shift from word-sequence evaluation toward *behavior-sequence* paradigms for longitudinal NLP.

2 Dataset

We use the PTSD-STOP dataset (Ringwald et al., 2025) comprising dense longitudinal monitoring of PTSD symptoms with multimodal daily measurements. Participants were recruited through Stony

Brook’s World Trade Center Health Program¹ and completed a daily protocol for up to 90 days. Each day, participants (i) recorded a brief video diary on a personal smart device and (ii) completed a self-report PTSD symptom questionnaire (PCL; Ruggero et al., 2021). Participants received \$5 per completed daily entry (video + survey), for up to \$450 across the study.

	Total Data	Cross-Sectional	Prospective	Cross-Sectional & Prospective
# Participants	238	48	190	48
# Documents	17,051	2,290	4,541	1,147
Avg Docs per Person	71.6	47.7	23.9	23.9

Table 1: **Dataset statistics overall and within each evaluation regime.** Matrix icons depict how instances are partitioned for the split of interest: yellow cells indicate training person-day instances and green cells indicate the corresponding test instances (gray denotes unused). Statistics corresponds to the green region of respective column.

Daily diary videos and transcripts. In the video diaries, participants described their daily experiences and emotions in response to prompts (e.g., “tell me about the worst part of your day,” “describe when you felt most happy today”)². We transcribe videos using the Whisper automatic speech recognition model (Radford et al., 2023) and treat each transcript as a document. The average diary length is 4.52 minutes (SD=2.12) and contains 646.6 words on average (SD=358.6).

Outcome variable. Our primary outcome is the participant’s daily PTSD symptom severity score from the PCL questionnaire, denoted y_{it} for person i on day t . In this dataset, the PCL score ranges from 1 (low severity) to 5 (high severity).

Analytic sample and coverage. After filtering to participants with sufficiently dense longitudinal coverage (at least 50% days with available observations for a maximum of 90 days) and aligning diaries with outcome, our final analytic sample contains 238 participants with 17,051 documents (mean=71.6 documents/person). Table 1 summarizes dataset counts overall and within the evaluation subsets used throughout the paper.

Participant characteristics. Participants were predominantly White (60%; 4% Black/African American, 6% other/multi-racial, 29% not recorded/refused) and male (84%), with an average

¹www.stonybrookmedicine.edu/WTC

²refer §A for the list of prompted questions

age of 59.87 years (range 42–79). We refer readers to Ringwald et al. (2025) for additional recruitment, protocol, and cohort details.

Tasks. We study two language-to-symptom prediction settings. For the longitudinal evaluation paradigm (§3), we consider a same-day *nowcasting* task: predicting daily PTSD severity $y_{i,t}$ from same-day language $x_{i,t}$. For the longitudinal modeling analyses (§4), we consider *one-day-ahead forecasting*: predicting next-day severity $y_{i,t+1}$ from prior-day language $x_{i,t}$.

Data preprocessing steps used for creating the datasets analyzed for the evaluation paradigm (§3) is detailed in Appendix §B.1 and that used for longitudinal modeling (§4) is detailed in Appendix §B.2.

Illustration of splits. For ease of reading, each table/figure is accompanied by a small matrix icon indicating the evaluation regime at the top. Rows correspond to people i and columns to days t ; each cell represents an observed instance (e.g., $(x_{i,t}, y_{i,t})$ for nowcasting, or $(x_{i,t}, y_{i,t+1})$ for forecasting). Yellow cells denote training instances, green cells denote test instances for the evaluation split of interest, and gray cells denote instances not used for that split. Cross-sectional evaluation splits by people ($i \in \mathcal{I}_{\text{train}}$ vs. $i \in \mathcal{I}_{\text{test}}$), prospective evaluation splits by time ($t \leq \tau$ vs. $t > \tau$), and cross-sectional & prospective splits by both (train: $i \in \mathcal{I}_{\text{train}}, t \leq \tau$; test: $i \in \mathcal{I}_{\text{test}}, t > \tau$).

Missingness handling. When language is missing on a given day, we impute $x_{i,t}$ by carrying forward the most recent available language observation for that participant (last observation carried forward). In contrast, we do not impute outcomes: days with missing PCL scores ($y_{i,t}$) are excluded from training and evaluation. Accordingly, the forecasting task uses $(x_{i,t}, y_{i,t+1})$ pairs only when $y_{i,t+1}$ is observed.

3 Evaluation: Splits and Metrics for Longitudinal NLP

Traditional evaluation yielded substantially different errors than ecologically valid settings.

To demonstrate the issues ensued by traditional paradigm, we mirror common settings where many documents come from a limited number of human sources (Geva et al., 2019) by sampling a subset of 20 individuals from PTSD-STOP³. Across evaluation regimes, the *same* document-level modeling

³details in Appendix §B.1

	Ecologically Improbable	Ecologically Valid	
		Cross Sectional	Prospective
Cross Validation Split:	Traditional Test Set mae	Test Set mae	Test Set mae
Typical model	.520	.757 [‡]	.455 [‡]
Baseline: Mean of Train	.660	.619	.598
N docs (train / test)	1008 / 421	1008 / 421	1008 / 421
Mean (Train / Test)	1.73 / 1.78	1.81 / 1.60	1.78 / 1.67
Std Dev (Train / Test)	0.79 / 0.84	0.84 / 0.71	0.84 / 0.72

Table 2: **Traditional evaluation suffers the ecological fallacy and provides no insight about generalization over people or time dimension.** This demonstration shows that ecologically invalid train - test split can lead to substantially different accuracies from valid settings using a typical model. This also shows that typical models do not necessarily work well for ecologically valid settings. For example, it leads to poor cross-sectional generalization as existing training methods regularize at document-level and not person-level. Data was sampled from 20 random people from PTSD-STOP data spanning a maximum of 90 days. *Typical model*: RoBERTa-large with a task regression head and L2 regularization. [‡] $p < .001$ (one-sided paired t-test vs. baseline).

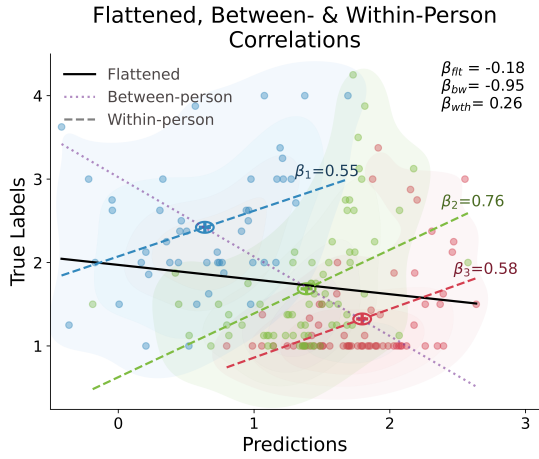
pipeline led to sharply different, and sometimes reversed conclusions (Table 2). Using a typical document-level model⁴ and a baseline that predicts the training-set mean, we find that under a *traditional* random document split the model outperformed the baseline (MAE: .520 vs .660; $\Delta\text{MAE} = -.140$). However, under a more realistic *cross-sectional* split (test on unseen people), this conclusion reversed: the model underperformed the baseline (MAE: .757 vs .619; $\Delta\text{MAE} = +.138$; $p < .001$). Under a *prospective* split (test on future days for the same people), the model again outperformed the baseline (MAE: .455 vs .598; $\Delta\text{MAE} = -.143$; $p < .001$), but the estimated error differed substantially from the traditional split (.455 vs .520).

These discrepancies arise because longitudinal NLP data are person-indexed and time-ordered: random document splits can leak person identity across train and test and can induce temporal leakage (training on later observations while evaluating on earlier ones). Such an interpolation-style evaluation is typically easier than deployments that require generalization to *new people* or *future days*.

Importantly, train/test outcome means and stan-

⁴Fine-tuning the task-specific/regression layer of RoBERTa-large (Liu et al., 2019). Refer to Appendix §C for details.

(a) Demonstration of Flattened metric obscuring cross-sectional and temporal evaluation of the model.



(b) Between- and Within-Person metrics capture how well the model generalizes over people and time.

Metric Scope	Cross-sectional Test Set		Prospective Test Set		Cross-sectional & Prospective	
	MAE (↓)	r(↑)	MAE (↓)	r(↑)	MAE (↓)	r(↑)
Flattened	.656	.295	.434	.655	.571	.314
Between-Person	.524	.450	.230	.962	.409	.468
Within-Person	.654	.285	.426	.297	.578	.214

Figure 1: **Flattened metrics can mask whether models learn people or dynamics.** *Left:* For three illustrative users from cross-sectional test set, predictions covary within individuals, but the between-person relationship of person-level means is reversed in direction; a pooled document-level fit conflates these effects, depicted by the flattened fit. *Right:* Decomposing performance into between-person and within-person MAE/r across split regimes separates individual differences from temporal variation and clarifies what drives apparent “good” performance.

standard deviations were similar across paradigms (Table 2), and the baseline remained relatively stable, suggesting that the observed differences primarily reflect the evaluation protocol rather than large distribution shifts. Overall, Table 2 shows that the *same* modeling pipeline can appear effective, ineffective, or even harmful depending on whether evaluation targets generalization across *documents*, *people*, or *time*.

METRIC SCOPE	EQUATION
FLATTENED	$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{t_i} f(y_{ij}, \hat{y}_{ij})$
BETWEEN-PERSON	$\frac{1}{n} \sum_{i=1}^n f(\bar{y}_i, \bar{\hat{y}}_i)$
WITHIN-PERSON	$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{t_i} \sum_{j=1}^{t_i} f(y_{ij}, \hat{y}_{ij}) \right)$

Table 3: **Metric scopes for person-indexed longitudinal evaluation.** Flattened metrics pool all person-days, while between-person and within-person metrics separately evaluate person-level differences and within-person temporal variation. Here f is the per-instance metric function, and y_{ij}, \hat{y}_{ij} are the observed and predicted outcomes for person i at time j .

A single pooled metric hid whether the model learned people or dynamics. Even with realistic splits, standard “flattened” metrics (computed by pooling all documents together) mixed two distinct

prediction goals: (i) capturing stable differences between people (who tended to have higher scores on average) and (ii) capturing within-person change over time (day-to-day deviations). These are both important in longitudinal behavioral prediction, but a model can perform well on one while performing poorly on the other. We therefore propose reporting complementary *between-person* and *within-person* metrics (Table 3). Between-person metrics evaluate predictions after aggregating within each person (e.g., comparing each person’s mean predicted score to their mean true score), while within-person metrics evaluate performance within each person first and then average across people. In simple terms, between-person evaluates how well the model captured *who was higher*, and within-person evaluates how well it captured *when someone was higher*.

Between- and within-person metrics reveal what drives “good” performance. Figure 1a illustrates why pooled evaluation could be misleading. For three example individuals from the cross-sectional test set, predictions covaried with the true labels within each person over time, yet the relationship between person-level averages differed in direction in this illustrative subset. When documents were pooled, these sources of variation were blended into a single summary that obscured which component the model captured.

Figure 1b makes this distinction explicit by

decomposing performance into between-person (individual differences) and within-person (temporal dynamics) components.⁵ Across all regimes, between-person performance exceeds within-person performance. On the *cross-sectional* test set, between-person error and correlation are better than within-person (MAE: .524 vs .654; r : .450 vs .285), indicating that the model more reliably captures stable person-level differences than day-to-day variation. The gap is larger on the *prospective* test set: between-person correlation is near-perfect ($r = .962$; MAE=.230) while within-person correlation remains modest ($r = .297$; MAE=.426), implying that the strong flattened prospective correlation (flattened $r = .655$) is driven primarily by person-level signal. The combined cross-sectional & prospective setting shows the same pattern (between-person $r = .468$ vs within-person $r = .214$). Overall, reporting between- and within-person metrics alongside pooled metrics provides a clearer account of model behavior in longitudinal NLP: whether performance is driven by learning *who* differs or by tracking *how* individuals change over time.

4 Modeling: From Isolated Documents to Behavioral Sequences

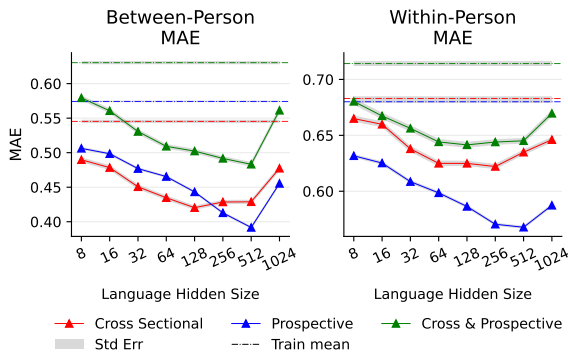


Figure 2: **Between- and within-person MAE vs. representation size.** Error shows a U-shaped trend across regimes: cross-sectional and cross-sectional & prospective perform best with small representations (64–128), whereas prospective benefits from larger size (~512).

Representation capacity depended on the generalization target (people vs. time). We first considered the typical NLP modeling setup that maps an *isolated* document representation to an

⁵We resample the data to include the third evaluation split; details are in Appendix §B.1.

outcome, and varied the dimensionality of the language representation (“hidden size”). Prior work in human-level NLP suggested that relatively low-dimensional representations can suffice for prediction (V Ganesan et al., 2021); applying our evaluation paradigm added an important nuance. We reduced the dimensions of language using Principal Component Analysis.

As shown in Figure 2, the optimal dimensionality shifted depending on whether evaluation required generalization to *new people* (cross-sectional), to *future days* (prospective), or both (cross-sectional & prospective). In the cross-sectional and cross-sectional & prospective test sets, performance was best at relatively low dimensionality (64-128), consistent with prior qualitative finding that increasing representation size does not monotonically improve human-level prediction. In contrast, for the prospective test set, error decreased further at higher dimensionality (around 512), and this pattern held for both between-person and within-person MAE.

Across settings, prospective errors were lower than cross-sectional errors, indicating that generalizing to *unseen people* was more difficult than generalizing forward in *time* for seen people. Finally, all language-based models outperformed the “mean of train” baseline (dotted lines), and the same qualitative trends were observed under SMAPE (Figure 5).

Modeling language as behavioral sequences improves prediction under ecologically valid evaluation. We next tested whether treating daily diaries as a *sequence of behaviors*, rather than independent documents, improves performance. Using an autoregressive ridge model, we varied the history length h (number of prior days of language representations provided as input), where $h=1$ corresponds to the traditional setup.

As shown in Figure 3, incorporating longer histories reduced error across evaluation sets, with the most consistent gains in the cross-sectional split, where MAE improved monotonically with h . In the prospective split, autoregressive performance improved initially and then plateaued as history increased, while in the cross-sectional & prospective split, autoregressive performance exhibited a U-shaped relationship with h , suggesting an accuracy – complexity trade-off under the most stringent generalization setting.

Across splits, the best-performing sequence mod-

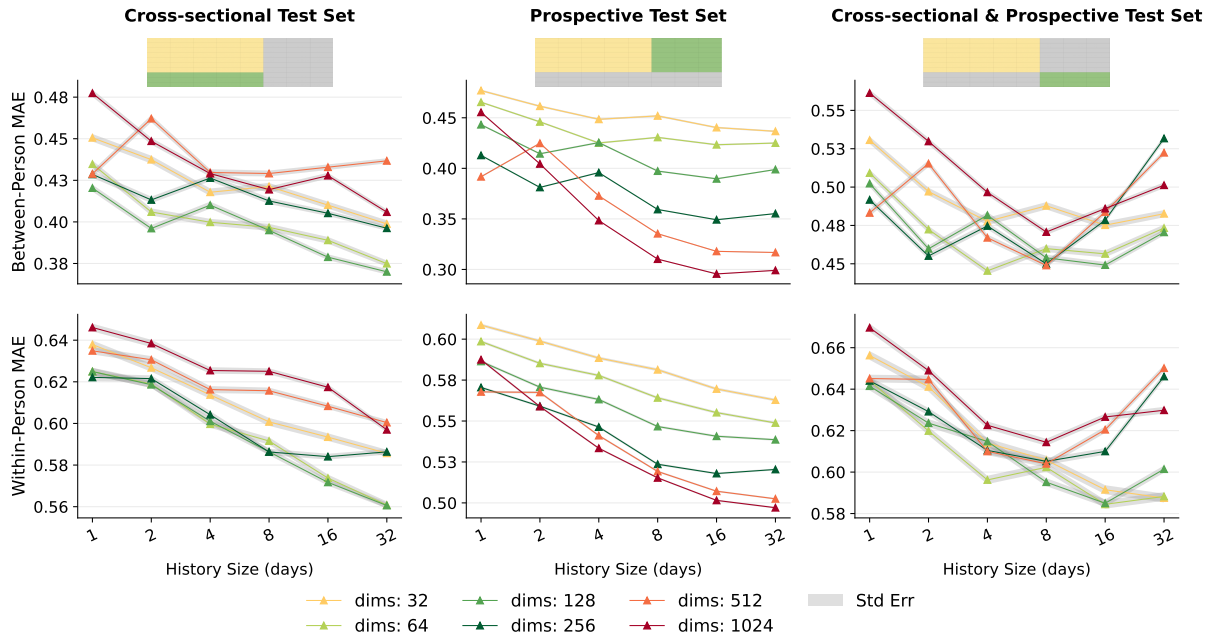


Figure 3: **Between- and within-person MAE vs. history length.** Using longer history generally improves performance, but the best history–capacity trade-off depends on the regime: prospective benefits from longer context with larger representations (≥ 512), while cross-sectional and cross-sectional & prospective improve primarily with longer context at smaller size (~ 64). Across splits, the average best size is 128.

els achieved improvements that exceeded the uncertainty in our MAE estimates (standard errors were non-overlapping at the best-performing settings), indicating that the gains are robust and not driven by noise. Together, these results motivate *sequence modeling of language as the default* for longitudinal behavioral prediction.

The optimal temporal inductive bias depends on the generalization target, implying a need for state coarseness. To understand what aspects of history modeling drive these gains, we compared three increasingly expressive ways to map a length- h language history into a predictive state: (i) a *bagged history* (BoE) that averages the last h days’ representations into a single coarse summary vector, (ii) an *autoregressive* model (AR) that explicitly parameterizes lagged dynamics over the h days, and (iii) a small *transformer* that can model content-dependent *interactions* across days (trained from scratch with a causal mask limiting attention to the past h days; Figure 4).

Strikingly, the best-performing inductive bias depended on the evaluation split. In the cross-sectional split, AR performed best, with the transformer close behind, suggesting that structured dynamics help while more flexible interaction modeling provides limited additional benefit. In the

prospective split, the transformer clearly outperformed AR and BoE, consistent with the presence of predictive signal in temporal interactions once individual baselines are observed. In the cross-sectional & prospective split, BoE performed best, indicating that a *coarser*, more regularized state is advantageous under the most stringent generalization setting.

Overall, no single temporal inductive bias dominates; instead, models benefit from an internal structure that can accommodate different *levels of state coarseness*, ranging from pooled summaries, to explicit dynamics, to interaction-rich sequence models, depending on whether the task requires generalization across people, time, or both.

5 Discussion

The independence assumption is (implicitly) a human-source assumption. Our results show that assuming documents are independent is rarely an inert or inconsequential abstraction in NLP; it implicitly assumes the *human sources* that generate supervision (authors, annotators) are exchangeable and independent. Because supervision comes from a finite set of people, models can exploit source-specific regularities and fail to generalize to unseen sources, underscoring the importance of documenting dataset provenance (Geva et al., 2019; Bender

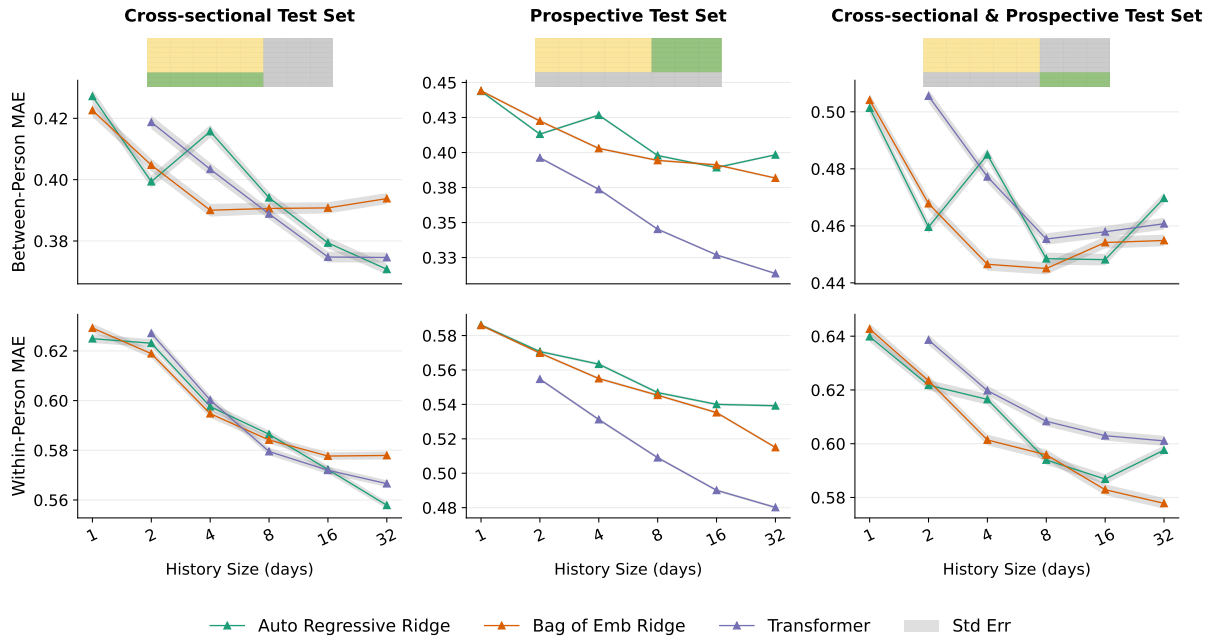


Figure 4: **AR vs. BoE vs. Transformer across history length (128 dims per day)**. Between-person (top) and within-person (bottom) MAE as a function of history h . AR performs best for cross-sectional generalization, while modeling temporal interactions (Transformer) yields the largest gains for prospective generalization but performs worst for cross-sectional & prospective.

and Friedman, 2018). This aligns with human-centered NLP that treats language as *verbal behavior* produced by people in context (Boyd and Schwartz, 2021; Soni et al., 2024). Our evaluation framework operationalizes this perspective by making the intended generalization target explicit: *new documents from seen people, new people, and/or future time*.

Model evaluations lose meaning and value when they lack veridicality Our cross-sectional and prospective splits instantiate evaluation under distribution shift: performance depends on whether test-time data share the same sources (people) and time period as training, and practical deployments often restrict this overlap (e.g., predicting for *unseen individuals* or *future days*). This framing aligns with in-the-wild robustness work and evidence that models can learn “shortcuts” that look strong under convenient test sets but break under shifted conditions (Koh et al., 2021; Geirhos et al., 2020; Rosenblatt et al., 2024). In longitudinal NLP, temporally aware evaluation has been motivated as necessary for understanding behavioral dynamics rather than static aggregates (Matero and Schwartz, 2020; Tsakalidis et al., 2022). Our findings reinforce this methodological point for NLP: *evaluation design is not merely an implementation choice; it changes the*

scientific claim a paper is able to support.

Metrics should reflect whether we care about people, time, or both. When documents are repeated measures of individuals, a single pooled (flattened) score can misalign with the intended behavioral claim by mixing stable between-person differences with within-person change (cf. ecological correlations Robinson, 1950). Reporting *between-person* and *within-person* variants of standard metrics makes these targets explicit and helps diagnose whether a model primarily recovers stable baselines versus tracks day-to-day variation. This mirrors measurement practice in psychology and longitudinal designs, where separating within- and between-person signal is foundational to interpretation (Shiffman et al., 2008). For NLP readers, the practical recommendation is lightweight: keep standard pooled metrics for comparability, but add target-aligned decompositions when the data are person-indexed over time.

Implications for modeling and dataset construction. Treating documents as behavioral sequences surfaces modeling choices that independence evaluation can hide, including conditioning on temporal context to model states and trajectories rather than isolated utterances (Matero et al., 2021; Soni et al., 2022). Although human-

centered shared tasks increasingly highlight longitudinal structure and temporally sensitive evaluation, these practices have not yet become default in mainstream NLP benchmarks (Tsakalidis et al., 2022; Matero and Schwartz, 2020). We therefore recommend a reporting norm for datasets with finite human sources or repeated measures: describe source structure, justify the generalization target, and align splits and metrics to that target (Bender and Friedman, 2018). Taken together, these steps move longitudinal NLP toward a methodology where conclusions are about *people and time*; not only about fitting collections of documents.

6 Related Work

Target-aligned evaluation for longitudinal NLP. A growing body of work recognizes that many applied NLP settings require temporally sensitive evaluation, especially in mental health and affective modeling (Matero and Schwartz, 2020; Tsakalidis et al., 2022). However, longitudinal evaluations often operationalize only one axis at a time (e.g., within-user temporal change or user-level holdout), making it hard to interpret what kind of generalization a reported score actually reflects. Our work complements this line by explicitly separating *cross-sectional*, *prospective*, and *cross-sectional & prospective* regimes, so evaluation claims are aligned with deployment targets.

Between-person vs. within-person variation. In “human-level” prediction, several studies emphasize that the effective sample size is the number of *people* rather than *documents*, and that performance is strongly shaped by person-level variation (V Ganesan et al., 2021). Conversely, longitudinal NLP work has highlighted within-person change as a core target (Tsakalidis et al., 2022). To our knowledge, prior work rarely unifies these two sources of variation into a single evaluation framework; our split regimes and between-/within-person metrics are designed to jointly diagnose generalization across people *and* across time.

Source-specific leakage goes beyond coarse human attributes. Related concerns appear in work on non-i.i.d. supervision, showing that models can exploit source-specific regularities and fail to generalize when the human source changes. For example, Geva et al. (2019) demonstrate that models can leverage annotator identifiers and that performance drops under annotator-disjoint splits, moti-

vating source-aware evaluation. More recently, Orlikowski et al. (2023) show that even when sociodemographic attributes correlate with label variation, modeling those attributes does not recover individual annotator behavior, highlighting that leakage cannot be reduced to coarse metadata; our setting instantiates an analogous issue for *authors* with the added complication of temporal ordering.

Modeling people over time. A parallel line of work adapts models to individuals via user representations or personalization, and human-language modeling arguments similarly treat individuals as generators of text with states and traits (Lynn et al., 2017; Soni et al., 2022). In longitudinal prediction tasks, autoregressive and sequence-aware models have been used to exploit temporal context (Matero and Schwartz, 2020). Our contribution is complementary: we show that making the evaluation target explicit changes which modeling choices look favorable, motivating sequence inputs by default and model internals that support different coarseness of state representations depending on the regime.

7 Conclusion

We argue that traditional NLP evaluation practice implicitly assumes independence between documents, even when datasets are generated by a finite set of humans and contain repeated measures over time. Using a temporally dense daily-diary dataset, we showed that this random document-level splits can yield substantially different, and sometimes reversed conclusions relative to ecologically valid evaluation targets. To make evaluation claims explicit, we operationalized split paradigms that separately test generalization to unseen people (cross-sectional) and to future days (prospective).

We further showed that standard pooled metrics can hide what a model actually learned in longitudinal settings, and that reporting between-person and within-person metrics clarifies whether performance is driven by stable individual differences or sensitivity to day-to-day change. Finally, we connected evaluation to modeling: once the target is specified, incorporating temporal context becomes a principled modeling choice, and its benefits (and failure modes) become measurable. Overall, our goal is not a benchmark race on a single endpoint, but a practical methodological shift toward longitudinal NLP as modeling *people and time*, with evaluation protocols that support the scientific and deployment claims we want to make.

575 **Limitations**

576 **Scope of methods and models studied.** This paper
577 is a methodology paper rather than a benchmark
578 race, and our experiments intentionally focus on
579 simple, transparent model classes (ridge regression
580 and autoregressive variants) to isolate evaluation
581 effects. As a result, we do not claim that the specific
582 optimal representation dimensionalities or history
583 lengths we observed will transfer unchanged
584 to other datasets, tasks, or model families. We
585 leave such line of work (V Ganesan et al., 2021;
586 Singh et al., 2025) to future works. More expressive
587 sequential architectures, additional modalities,
588 or alternative training objectives may change absolute
589 performance (Matero et al., 2021; Rao et al.,
590 2025), but the core risk we highlight—that convenient
591 independence-assumed evaluation can reward
592 reliance on source-specific regularities and fail under
593 shifted regimes—has been observed broadly in
594 ML under distribution shift and shortcut learning
595 (Koh et al., 2021; Geirhos et al., 2020). Finally,
596 we study one clinical endpoint (daily PCL) as a
597 concrete case; extending this framework to other
598 outcomes (e.g., functional impairment, treatment
599 response) and to other longitudinal NLP settings is
600 an important direction for future work.

601 **Dataset- and measurement-specific constraints.**

602 Our empirical demonstrations are grounded in a
603 single intensive longitudinal cohort (PTSD-STOP)
604 with daily self-reports and diary entries, which
605 enables the methodological analyses we pursue
606 but also constrains external validity. First, the cohort’s
607 demographics and recruitment context (a clinically
608 monitored trauma-exposed population) may not represent
609 other clinical groups or the general population, and
610 prior work has cautioned that clinical NLP systems
611 can inherit and amplify such representational skews
612 if generalization targets and subgroup performance
613 are not made explicit (Bear Don’t Walk et al., 2022;
614 Obermeyer et al., 2019). Second, intensive longitudinal
615 data are rarely missing-at-random: adherence varies
616 over time and across individuals, and missingness
617 mechanisms can bias both modeling and evaluation
618 if not carefully characterized (Shiffman et al., 2008;
619 Stone et al., 2023; Heron et al., 2017). Third, we
620 rely on automatic speech recognition to obtain transcripts;
621 transcription errors and diarization artifacts can
622 introduce additional noise that may differentially
623 affect within-person versus between-person signals,
624 and should be treated as part of the mea-

625 surement process rather than ignored (Shiffman
626 et al., 2008; Stone et al., 2023).

627 **Ethical Considerations**

629 **Sensitive human-subject data and governance.**

630 This work analyzes longitudinal materials linked
631 to mental health symptom reports, which are inherently
632 sensitive and may contain personally identifying
633 information (even after transcription). Accordingly,
634 analyses should be conducted under appropriate
635 human-subject oversight and data-use agreements
636 consistent with established principles for research
637 with human participants (U.S. Department of Health
638 & Human Services, 1979). Because de-identification
639 is not a guarantee of anonymity, especially for rich
640 narrative data, we follow the conservative stance
641 that access controls and minimization of shared
642 artifacts (e.g., avoiding release of raw text/video)
643 are often necessary for participant protection
644 (U.S. Department of Health & Human Services,
645 2025; Ayers et al., 2018). We also view dataset
646 provenance and documentation as part of ethical
647 reporting in NLP: clearly stating who produced the
648 data (sources), under what conditions, and what
649 the intended generalization target is helps prevent
650 misuse and misinterpretation of results (Bender
651 and Friedman, 2018; Gebru et al., 2021).

653 **Risks of downstream use in mental health contexts.**

654 Models that predict mental health-related
655 outcomes can be misapplied in ways that harm
656 individuals: erroneous inferences may contribute
657 to stigma, inappropriate monitoring, or decisions
658 made without clinical context. Prior work in
659 digital mental health NLP emphasizes that prediction
660 should not be conflated with diagnosis, that
661 construct validity is difficult to establish from
662 language traces alone, and that failure modes are
663 often context-dependent (Chancellor and De
664 Choudhury, 2020; Ernala et al., 2019). More
665 broadly, the field has called for explicit ethics
666 disclosures and clearer articulation of intended
667 use, beneficiaries, and plausible harms in
668 mental health prediction research (Ajmani et al.,
669 2023; Chancellor et al., 2019). Consistent with
670 this, our framing and experiments are aimed at
671 improving *methodology* (what we can claim from
672 evaluation), not at promoting clinical deployment.

673 **Re-identification and leakage considerations.**

674 Longitudinal diaries are especially vulnerable to

675	re-identification because they accumulate unique	Services, 2025; Ayers et al., 2018).	726
676	life details over time. Even when direct identifiers		
677	are removed, verbatim excerpts can enable reverse	Use of AI assistance. We used AI assistants to	727
678	identification, and model artifacts can inadvertently	support drafting (e.g., paraphrasing for clarity) and	728
679	memorize or surface sensitive details (Ayers et al.,	coding/formatting during manuscript preparation.	729
680	2018; Chancellor et al., 2019). For this reason, re-	All generated text and code were reviewed and ver-	730
681	leasing raw text, audio, or video should be treated	ified by the authors, and all co-authors participated	731
682	as a high-risk action; safer alternatives include re-	in editing and auditing to ensure accuracy, appro-	732
683	leasing code, evaluation splits, and aggregate statis-	appropriate attribution, and consistency with the study	733
684	tics, or providing controlled access mechanisms	protocol and results.	734
685	where appropriate (U.S. Department of Health &		
686	Human Services, 2025; Gebru et al., 2021).		
687	Fairness, subgroup validity, and clinical safety.	References	735
688	Clinical NLP systems can encode uneven perfor-	Leah Hope Ajmani, Stevie Chancellor, Bijal Mehta,	736
689	mance across demographic or clinical subgroups,	Casey Fiesler, Michael Zimmer, and Munmun	737
690	potentially exacerbating existing inequities if used	De Choudhury. 2023. <i>A systematic review of ethics</i>	738
691	in practice (Bear Don't Walk et al., 2022; Ober-	<i>disclosures in predictive mental health research.</i> In	739
692	meyer et al., 2019). A practical implication for	<i>Proceedings of the 2023 ACM Conference on Fair-</i>	740
693	longitudinal NLP is that “ecologically valid” splits	<i>ness, Accountability, and Transparency (FAccT),</i>	741
694	should be complemented with subgroup-aware re-	pages 1311–1323.	742
695	porting whenever sample sizes permit, and that	John W. Ayers, Theodore L. Caputi, Camille Nebeker,	743
696	claims should be scoped to the population repre-	and Mark Dredze. 2018. <i>Don't quote me: reverse</i>	744
697	sented by the data (Mitchell et al., 2019; Bender	<i>identification of research participants in social media</i>	745
698	and Friedman, 2018). Finally, mental health appli-	<i>studies.</i> <i>npj Digital Medicine</i> , 1:30.	746
699	cations require particular caution: as emphasized	Oliver J. Bear Don't Walk, H. Reyes Nieva, S. S. J.	747
700	by the CLPsych community, ethical deployment	Lee, and Noémie Elhadad. 2022. <i>A scoping review</i>	748
701	demands careful consideration of consent, stake-	<i>of ethics considerations in clinical natural language</i>	749
702	holder involvement, and clinical safety boundaries	<i>processing.</i> <i>JAMIA Open</i> , 5(2):ooac039.	750
703	beyond standard ML evaluation (Orr et al., 2022).	Emily M. Bender and Batya Friedman. 2018. <i>Data</i>	751
704	Sensitive human-subject data and governance.	<i>statements for natural language processing: Toward</i>	752
705	This work analyzes longitudinal materials linked	<i>mitigating system bias and enabling better science.</i>	753
706	to mental health symptom reports, which are inher-	<i>Transactions of the Association for Computational</i>	754
707	ently sensitive and may contain personally identify-	<i>Linguistics</i> , 6:587–604.	755
708	ing information (even after transcription). Accord-	Ryan L. Boyd and H. Andrew Schwartz. 2021. <i>Natu-</i>	756
709	ingly, analyses should be conducted under appro-	<i>ral language analysis and the psychology of verbal</i>	757
710	appropriate human-subject oversight and data-use agree-	<i>behavior: The past, present, and future states of the</i>	758
711	ments consistent with established principles for re-	<i>field.</i> <i>Journal of Language and Social Psychology</i> ,	759
712	search with human participants (U.S. Department	40(1):21–41.	760
713	of Health & Human Services, 1979). The study pro-	Stevie Chancellor, Michael L. Birnbaum, Eric D. Caine,	761
714	cedure was reviewed and approved by the relevant In-	Vincent M. Silenzio, and Munmun De Choudhury.	762
715	stitutional Review Board (IRB), and all researchers	2019. <i>A taxonomy of ethical tensions in inferring</i>	763
716	followed institutional Human Subjects Research	<i>mental health states from social media.</i> In <i>Proceed-</i>	764
717	guidelines for secure data handling (e.g., restricted	<i>ings of the 2019 ACM Conference on Fairness, Ac-</i>	765
718	access, secured storage/compute, and minimization	<i>countability, and Transparency (FAccT),</i> pages 79–	766
719	of data movement and shared artifacts). Because	88.	767
720	de-identification is not a guarantee of anonymity,	Stevie Chancellor and Munmun De Choudhury. 2020.	768
721	especially for rich narrative data, we follow the con-	<i>Methods in predictive techniques for mental health</i>	769
722	servative stance that access controls and minimiza-	<i>status on social media: a critical review.</i> <i>npj Digital</i>	770
723	tion of shared artifacts (e.g., avoiding release of	<i>Medicine</i> , 3(1):43.	771
724	raw text/video) are often necessary for participant	Minje Choi, Jiaxin Pei, Sagar Kumar, Chang Shu, and	772
725	protection (U.S. Department of Health & Human	David Jurgens. 2023. <i>Do LLMs understand social</i>	773
		<i>knowledge? evaluating the sociability of large lan-</i>	774
		<i>guage models with SocKET benchmark.</i> In <i>Proceed-</i>	775
		<i>ings of the 2023 Conference on Empirical Methods in</i>	776
		<i>Natural Language Processing</i> , pages 11370–11403,	777
		Singapore. Association for Computational Linguis-	778
		tics.	779

780	Patrick J Curran and Daniel J Bauer. 2011. The disaggregation of within-person and between-person effects in longitudinal models of change. <i>Annual review of psychology</i> , 62(1):583–619.	837
781		838
782		839
783		840
784	Sindhu K. Ernala, Michael L. Birnbaum, Asra F. Rizvi, William A. Sterling, and Munmun De Choudhury. 2019. Methodological gaps in predicting mental health states from social media: Triangulating diagnostic signals. In <i>Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI)</i> , pages 1–16.	841
785		842
786		843
787		844
788		845
789		846
790		847
791	Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. Datasheets for datasets. In <i>Communications of the ACM</i> , volume 64, pages 86–92.	848
792		849
793		850
794		851
795		852
796	Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A. Wichmann. 2020. Shortcut learning in deep neural networks. <i>Nature Machine Intelligence</i> , 2:665–673.	853
797		854
798		855
799		856
800		857
801	Mor Geva, Yoav Goldberg, and Jonathan Berant. 2019. Are we modeling the task or the annotator? an investigation of annotator bias in natural language understanding datasets. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1161–1166, Hong Kong, China. Association for Computational Linguistics.	858
802		859
803		860
804		861
805		862
806		863
807		864
808		865
809		866
810	Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 708–719, New Orleans, Louisiana. Association for Computational Linguistics.	867
811		868
812		869
813		870
814		871
815		872
816		873
817		874
818	Ellen L Hamaker, Rebecca M Kuiper, and Raoul PPP Grasman. 2015. A critique of the cross-lagged panel model. <i>Psychological methods</i> , 20(1):102.	875
819		876
820		877
821	Adi Haviv, Ori Ram, Ofir Press, Peter Izsak, and Omer Levy. 2022. Transformer language models without positional encodings still learn positional information. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 1382–1390, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	878
822		879
823		880
824		881
825		882
826		883
827		884
828	Kristin E. Heron, R. Scott Everhart, Susan M. McHale, and Joshua M. Smyth. 2017. Ecological momentary assessment (ema) of youth behavior: A systematic review and recommendations. <i>Journal of Pediatric Psychology</i> , 42(10):1087–1107.	885
829		886
830		887
831		888
832		889
833	Lesia Hoffman and Robert S Stawski. 2009. Persons as contexts: Evaluating between-person and within-person effects in longitudinal analysis. <i>Research in human development</i> , 6(2-3):97–120.	890
834		891
835		892
836		893
	Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva Reddy. 2023. The impact of positional encoding on length generalization in transformers. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23</i> , Red Hook, NY, USA. Curran Associates Inc.	
	Oscar NE Kjell, Sverker Sikström, Katarina Kjell, and H Andrew Schwartz. 2022. Natural language analyzed with ai-based transformers predict traditional subjective well-being measures approaching the theoretical upper limits in accuracy. <i>Scientific reports</i> , 12(1):3918.	
	Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Anand Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Li, Irena Jia, et al. 2021. WILDS: A benchmark of in-the-wild distribution shifts. In <i>Proceedings of the 38th International Conference on Machine Learning (ICML)</i> .	
	Sumeet Kumar and Kathleen Carley. 2019. Tree LSTMs with convolution units to predict stance and rumor veracity in social media conversations. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5047–5058, Florence, Italy. Association for Computational Linguistics.	
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. <i>CoRR</i> , abs/1907.11692.	
	Veronica Lynn, Youngseo Son, Vivek Kulkarni, Niranjan Balasubramanian, and H. Andrew Schwartz. 2017. Human centered NLP with user-factor adaptation. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 1146–1155, Copenhagen, Denmark. Association for Computational Linguistics.	
	Matthew Matero and H. Andrew Schwartz. 2020. Autoregressive affective language forecasting: A self-supervised task. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 2913–2923, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
	Matthew Matero, Nikita Soni, Niranjan Balasubramanian, and H. Andrew Schwartz. 2021. MeLT: Message-level transformer with masked document representations as pre-training for stance detection. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2959–2966, Punta Cana, Dominican Republic. Association for Computational Linguistics.	
	Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In	

1008	<i>Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4647–4660, Dublin, Ireland.	1064
1009	Association for Computational Linguistics.	1065
1010		1066
1011	U.S. Department of Health & Human Services.	1067
1012	1979. The belmont report: Ethical principles	1068
1013	and guidelines for the protection of human sub-	1069
1014	jects of research. https://www.hhs.gov/ohrp/	1070
1015	regulations-and-policy/belmont-report/ .	
1016	Last reviewed Aug 26, 2024.	
1017	U.S. Department of Health & Human Ser-	
1018	vices. 2025. Guidance regarding methods	
1019	for de-identification of protected health in-	
1020	formation in accordance with the HIPAA	
1021	privacy rule. https://www.hhs.gov/hipaa/	
1022	for-professionals/privacy/special-topics/	
1023	de-identification/index.html . Last reviewed	
1024	Feb 3, 2025.	
1025	Adithya V Ganesan, Yash Kumar Lal, August Nilsson,	
1026	and H. Andrew Schwartz. 2023. <i>Systematic evaluation</i>	
1027	<i>of GPT-3 for zero-shot personality estimation</i> . In	
1028	<i>Proceedings of the 13th Workshop on Computational</i>	
1029	<i>Approaches to Subjectivity, Sentiment, & Social Me-</i>	
1030	<i>dia Analysis</i> , pages 390–400, Toronto, Canada. Asso-	
1031	ciation for Computational Linguistics.	
1032	Adithya V Ganesan, Siddharth Mangalik, Vasudha	
1033	Varadarajan, Nikita Soni, Swanie Juhng, João Sedoc,	
1034	H. Andrew Schwartz, Salvatore Giorgi, and Ryan L	
1035	Boyd. 2024. <i>From text to context: Contextualizing</i>	
1036	<i>language with humans, groups, and communities for</i>	
1037	<i>socially aware NLP</i> . In <i>Proceedings of the 2024 Con-</i>	
1038	<i>ference of the North American Chapter of the Asso-</i>	
1039	<i>ciation for Computational Linguistics: Human Lan-</i>	
1040	<i>guage Technologies (Volume 5: Tutorial Abstracts)</i> ,	
1041	pages 26–33, Mexico City, Mexico. Association for	
1042	Computational Linguistics.	
1043	Adithya V Ganesan, Matthew Matero, Aravind Reddy	
1044	Ravula, Huy Vu, and H. Andrew Schwartz. 2021.	
1045	<i>Empirical evaluation of pre-trained transformers for</i>	
1046	<i>human-level NLP: The role of sample size and dimen-</i>	
1047	<i>sionality</i> . In <i>Proceedings of the 2021 Conference of</i>	
1048	<i>the North American Chapter of the Association for</i>	
1049	<i>Computational Linguistics: Human Language Tech-</i>	
1050	<i>nologies</i> , pages 4515–4532, Online. Association for	
1051	Computational Linguistics.	
1052	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	
1053	Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz	
1054	Kaiser, and Illia Polosukhin. 2017. Attention is all	
1055	you need. In <i>Proceedings of the 31st International</i>	
1056	<i>Conference on Neural Information Processing Sys-</i>	
1057	<i>tems, NIPS’17</i> , page 6000–6010, Red Hook, NY,	
1058	USA. Curran Associates Inc.	
1059	A PTSD-STOP Daily Diary Study	
1060	Open Ended Questionnaire Each day, partici-	
1061	pants completed a brief video diary consisting of 13	
1062	open-ended prompts about their day. The camera	
1063	began recording automatically and the recording	
	uploaded automatically upon submission. Partici-	1064
	pants were instructed to answer each prompt out	1065
	loud while facing the camera, to expand on their	1066
	responses (rather than giving short answers), and to	1067
	avoid reading the questions aloud. The submit but-	1068
	ton appeared after 3 minutes, with up to 10 minutes	1069
	available per prompt. The 13 prompts were:	1070
	1. Tell me about the best part of your day.	1071
	2. Tell me about the worst part of your day.	1072
	3. Describe when you felt most sad today.	1073
	4. Describe when you felt most scared or ner-	1074
	vous today.	1075
	5. Describe when you felt most annoyed today.	1076
	6. Describe when you felt most happy today.	1077
	7. How did you feel physically today? Did you	1078
	have any pain, discomfort, or other physical	1079
	symptoms? Please elaborate.	1080
	8. How did you get along with others today?	1081
	Please elaborate.	1082
	9. Did you have any unwanted, disturbing mem-	1083
	ories of a past stressful experience? Tell me	1084
	about this.	1085
	10. Today, did you avoid anything because it	1086
	would have made you uncomfortable? Tell	1087
	me more about what you did and why.	1088
	11. Did you feel on guard today? What made you	1089
	feel this way?	1090
	12. Describe anything that cheered you up today.	1091
	How did it go?	1092
	13. What happened today that you can feel thank-	1093
	ful for? Tell me about this.	1094
	Rating Scales In addition to the open-ended di-	1095
	ary, participants completed daily self-report rating	1096
	scales capturing (i) post-stressor/PTSD-related ex-	1097
	periences and overall stress that day, and (ii) expo-	1098
	sure to common daily stressors. Rating scales were	1099
	administered in two parts.	1100
	Part 1 (Symptom and stress ratings). Partici-	1101
	pants rated the following items everyday, keeping	1102
	their most stressful event in mind. Response op-	1103
	tions were: <i>Not at all, A little bit, Moderately, Quite</i>	1104
	<i>a bit, Extremely, and Skip.</i>	1105

- 1106 1. Today, I had repeated, disturbing, and un-
1107 wanted memories of the stressful experience.
- 1108 2. Today, I felt very upset because something
1109 reminded me of the stressful experience.
- 1110 3. Today, I avoided memories, thoughts, or feel-
1111 ings related to the stressful experience.
- 1112 4. Today, I avoided external reminders of the
1113 stressful experience (for example, people,
1114 places, conversations, activities, objects, or
1115 situations).
- 1116 5. Today, I felt distant or cut off from other peo-
1117 ple.
- 1118 6. Today, I had strong negative feelings such as
1119 fear, horror, anger, guilt, or shame.
- 1120 7. Today, I felt jumpy or easily startled.
- 1121 8. Today, I was “superalert” or watchful or on
1122 guard.
- 1123 9. Overall, how stressed did you feel today?
- 1124 **Part 2 (Daily stressor checklist).** Participants
1125 then indicated which of the following troublesome
1126 or stressful events occurred everyday (check all
1127 that apply):
- 1128 1. Had tension or argument with spouse, partner
1129 or close family
- 1130 2. Had tension or argument with others (e.g., co-
1131 worker, friend, etc)
- 1132 3. A lot of demands at home
- 1133 4. A lot of demands at job
- 1134 5. A lot of demands made by family
- 1135 6. Caring for a sick family member
- 1136 7. Problems with transportation
- 1137 8. Financial or money problem(s)
- 1138 9. Health-related event(s)
- 1139 10. Other troublesome things happened to me
- 1140 11. No troublesome or stressful things happened
1141 to me today

B Data Pre-processing 1142

B.1 Longitudinal Evaluation 1143

The violations of independence assumption arise whenever many documents are produced or labeled by a limited set of humans, and they are amplified in longitudinal data where documents are also ordered in time. Such dependencies are more often the rule than the exception. Even when an NLP task may appear on its surface to be “about the text itself” (e.g., classification, translation, summarization), supervision originates from a finite set of human sources: annotators create labels, and translators/summarizers create references. In practice, the “human set” can be surprisingly small relative to the number of documents, creating a deeply shared variance structure across instances. 1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157

For example, crowdsourced datasets often show extreme annotator imbalances: in the NLU datasets analyzed by Geva et al. (2019), MNLI contains 402k examples from 380 annotators and OpenBookQA contains 5k examples from 84 annotators, with the most prolific OpenBookQA annotator contributing 24% of all examples. Likewise, widely used summarization corpora contain reference summaries written within a limited set of institutions: Newsroom provides 1.3M article – summary pairs written by authors and editors across 38 news publications (Grusky et al., 2018). 1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169

Demonstration of ecological fallacy in traditional evaluation. Analyses for Table 2 and Figure 1a use the dataset constructed as follows. To mirror common settings where many documents come from a limited number of human sources (Geva et al., 2019), we form a controlled subset of 20 individuals from PTSD-STOP. 1170
1171
1172
1173
1174
1175
1176

We first filtered to individuals with non-trivial outcome variation, requiring at least 0.01 standard deviation in PCL over the full 90-day period, as well as within the first 60 days and last 30 days. We then sorted individuals by their mean PCL (averaged over days), binned them into 10 strata, and sampled two individuals uniformly at random from each stratum to obtain a balanced cohort of 20 people. 1177
1178
1179
1180
1181
1182
1183
1184
1185

We split this cohort into train/test in three ways: (a) a document-level random split (70/30), which is ecologically implausible for person-indexed data; (b) a cross-sectional split (70/30) by person, which holds out individuals; and (c) a prospective split (70/30) by time, using a temporal cutoff at day 1186
1187
1188
1189
1190
1191

1192	$\tau = 63$ (70% of 90 days). To match train/test sizes	1024 was used. The ridge penalty (L2 regulariza-	1241
1193	across these three settings, we randomly masked a	tion strength) is selected from $\{10^i : i \in [-2, 5]\}$	1242
1194	small number of person-day instances.	based on performance on a development set. To	1243
1195	Between- and within-person metrics isolate	avoid leakage, the development set is split from the	1244
1196	what drives performance. Analyses for Fig-	training data using the same evaluation regime as	1245
1197	ure 1b use the same 20-person cohort as above,	the corresponding test set (cross-sectional, prospec-	1246
1198	but evaluate one fixed model across multiple test	tive, or cross-sectional & prospective).	1247
1199	regimes. We set the prospective temporal cutoff		
1200	to a 67/33 split (i.e., $\tau \approx 60$) to define the train-	C.2 Longitudinal Modeling	1248
1201	ing window, and we additionally construct a cross-	Autoregressive ridge models. Autoregressive	1249
1202	sectional & prospective test set that requires gener-	ridge models take as input a length- h history of	1250
1203	alization over both people and time.	daily language representations and predict the tar-	1251
1204	Concretely, we train a single model on the	get outcome (nowcasting or one-day-ahead fore-	1252
1205	training region (yellow cells in the split illus-	casting, depending on the experiment). These mod-	1253
1206	tration) and evaluate that same model under	els are parameterized by history length h and per-	1254
1207	three regimes: cross-sectional (held-out people),	day representation size d , yielding an input dimen-	1255
1208	prospective (held-out future days), and cross-	sion of $h \cdot d$. Training and hyperparameter selec-	1256
1209	sectional & prospective (held-out people at held-	tion follow the same procedure as above (ridge	1257
1210	out future days). This design ensures differences	regression with the L2 penalty chosen on a regime-	1258
1211	in performance reflect the evaluation regime rather	matched development split).	1259
1212	than training different models.		
1213	B.2 Longitudinal Modeling	Bag-of-embeddings models. Bag-of-	1260
1214	For all results in §4, we use the full analytic sample	embeddings (BoE) models aggregate a length- h	1261
1215	of 238 participants (Table 1). For cross-sectional	history by averaging the h daily representations	1262
1216	evaluation, participants are stratified into train/test	into a single vector, which is then used for ridge	1263
1217	(80/20) based on their mean PCL over the study	regression. As a result, BoE is parameterized by	1264
1218	period. For prospective evaluation, we use a fixed	history length h and per-day representation size	1265
1219	temporal cutoff at day $\tau = 60$ for training, with	d , but its number of learned parameters scales	1266
1220	later days reserved for testing.	only with d (since the history is pooled before	1267
1221	C Model Training	prediction). We train BoE models with the same	1268
1222	C.1 Longitudinal Evaluation	ridge procedure and select the L2 penalty on a	1269
1223	Typical model. Our <i>typical</i> model is fine-tuning	regime-matched development split.	1270
1224	of the task-specific ridge regression layer of a	Transformers. Model architecture. To test	1271
1225	RoBERTa-large encoder (Liu et al., 2019). We	whether modeling <i>temporal interactions</i> improves	1272
1226	use this term because fine-tuned encoder mod-	prediction beyond pooled summaries (BoE) and	1273
1227	els remain a strong and widely used baseline for	linear dynamics (AR), we trained a minimal autore-	1274
1228	psychological and behavioral prediction from lan-	gressive transformer over the history sequence. We	1275
1229	guage (Kjell et al., 2022), and recent evidence	project each per-day representation from $d=128$	1276
1230	suggests that instruction-tuned LLMs do not con-	down to $d'=32$ using a learned linear map, and	1277
1231	sistently outperform fine-tuning smaller encoder	feed the resulting sequence into a transformer en-	1278
1232	models for behavioral prediction and psychological	coder with 1 layer and 1 attention head (hidden size	1279
1233	measurement (Singh et al., 2025; V Ganesan	32).	1280
1234	et al., 2023; Choi et al., 2023).	For comparability with AR and BoE, we made	1281
1235	In practice, we train ridge regression mod-	two design choices: (1) <i>No positional embeddings</i> .	1282
1236	els on transcript representations extracted from	We omit positional embeddings to reduce explicit	1283
1237	the second-to-last layer of RoBERTa-large using	order information, isolating gains attributable to in-	1284
1238	DLATK (Schwartz et al., 2017). Language dimen-	teraction modeling rather than architectural encod-	1285
1239	sions were reduced using Principal Component	ings of dynamics; prior work also suggests trans-	1286
1240	Analysis whenever a hidden dimension size of le	formers can recover positional information even	1287
		without explicit embeddings (Haviv et al., 2022;	1288
		Kazemnejad et al., 2023). (2) <i>Fixed history mask-</i>	1289
		<i>ing</i> . For a given history length h , we apply a lower-	1290

triangular attention mask and additionally prevent attention to tokens more than h days in the past by setting the corresponding attention logits to $-\infty$ prior to softmax. This ensures each transformer configuration is trained and evaluated under the same effective history constraint as the corresponding AR/BoE setting.

Hyperparameter optimization. We tuned the projection size d' , L2 weight decay, attention dropout, output-layer dropout, and learning rate on a regime-matched development split. We searched projection sizes among powers of two below 128 and selected $d'=32$. The selected hyperparameters were: weight decay = 1, attention dropout = 0.3, output dropout = 0.1, and learning rate = 10^{-3} with AdamW.

Additional checks. Although evaluating full-capacity transformers (e.g., larger depth/width, unmasked attention, or explicit sinusoidal/rotary positional embeddings (Vaswani et al., 2017; Su et al., 2024)) is beyond our scope, we verified that adding positional embeddings did not yield consistent improvements over our minimal configuration, and that increasing the maximum training history produced the same qualitative patterns reported in Figure 4.

D Metric Functions

We optimize models using the standard (flattened) mean squared error, leaving alternative objectives and metric formulations for future work. For evaluation, we report between- and within-person variants of Mean Absolute Error (MAE), Symmetric Mean Absolute Percentage Error (SMAPE), and Pearson correlation (r). SMAPE and r are bounded and scale-invariant, ranging in $[0, 2]$ and $[-1, 1]$, respectively.

SMAPE. For predictions \hat{y} and targets y , we compute

$$\text{SMAPE}(y, \hat{y}) = \frac{2}{N} \sum_{n=1}^N \frac{|\hat{y}_n - y_n|}{|y_n| + |\hat{y}_n| + \epsilon}, \quad (1)$$

where N is the number of evaluated instances and ϵ is a small constant to avoid division by zero. Since the outcomes range between $[1, 5]$, ϵ was set to 0.

Pearson correlation. We compute Pearson correlation between y and \hat{y} as

$$r(y, \hat{y}) = \frac{\sum_{n=1}^N (y_n - \bar{y})(\hat{y}_n - \bar{\hat{y}})}{\sqrt{\sum_{n=1}^N (y_n - \bar{y})^2} \sqrt{\sum_{n=1}^N (\hat{y}_n - \bar{\hat{y}})^2}}, \quad (2)$$

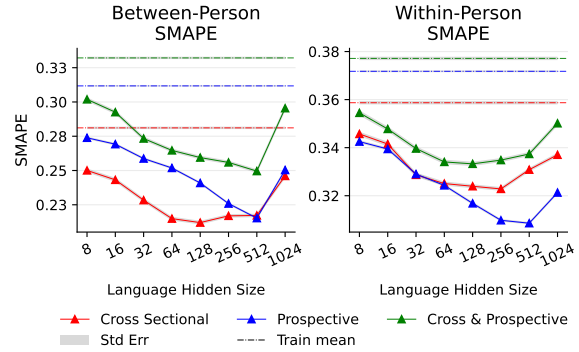


Figure 5: **Between- (top) and Within-Person (bottom) SMAPE as a function of hidden dimension size.** Forecasting performance follows a U-shaped trend as a function of hidden dimension size of language across all three evaluation sets. While a typical model requires only 64 dimensions of language for best performance on Cross-sectional and Cross-sectional & Prospective test sets, it requires 512 dimensions in Prospective evaluation set. Based on the best performance achieved in different settings, generalization to Cross-sectional & Prospective is the hardest, followed by Cross-sectional set and finally prospective set. Generalization to unseen people is harder than unseen time and within-person changes is harder than between-person differences.

where $\bar{y} = \frac{1}{N} \sum_{n=1}^N y_n$ and $\bar{\hat{y}} = \frac{1}{N} \sum_{n=1}^N \hat{y}_n$.

1336

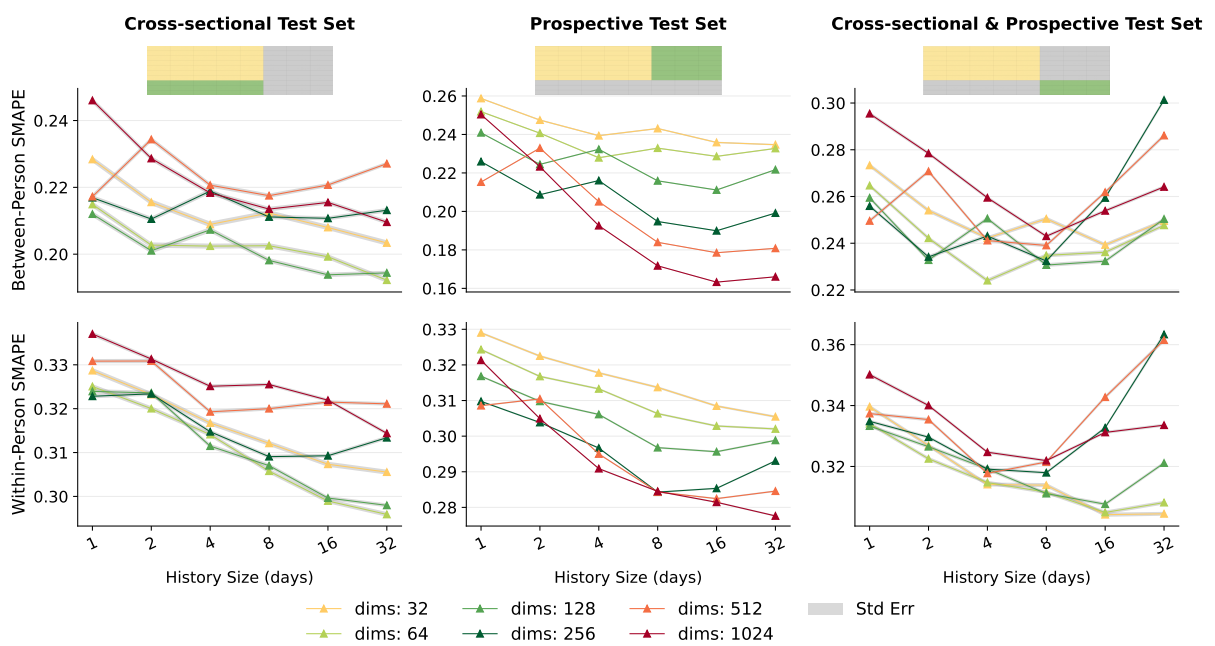
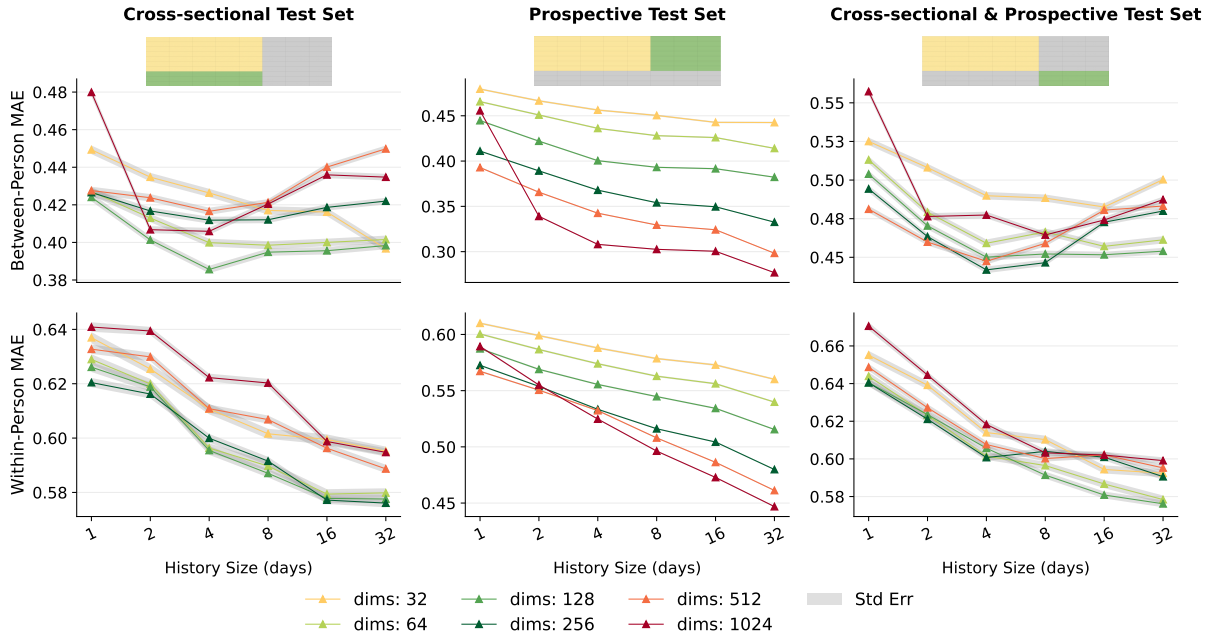


Figure 6: **Between- and Within-Person SMAPE as a function of history length for Auto Regressive Model.** Predictive Performance improves with modeling the temporal dynamics of linguistic behavior. For Prospective test set, between- and within-person performance improves with longer temporal context and higher dimensional sizes (512 and 1024). For Cross-sectional and Cross-sectional & Prospective evaluation sets, performance improves with temporal context at lower dimensions (hidden size=64).

(a) Between- (*top*) and Within-Person (*bottom*) MAE as a function of history length for Bag of Embeddings Model.



(b) Between- (*top*) and Within-Person (*bottom*) SMAPE as a function of history length for Bag of Embeddings Model.

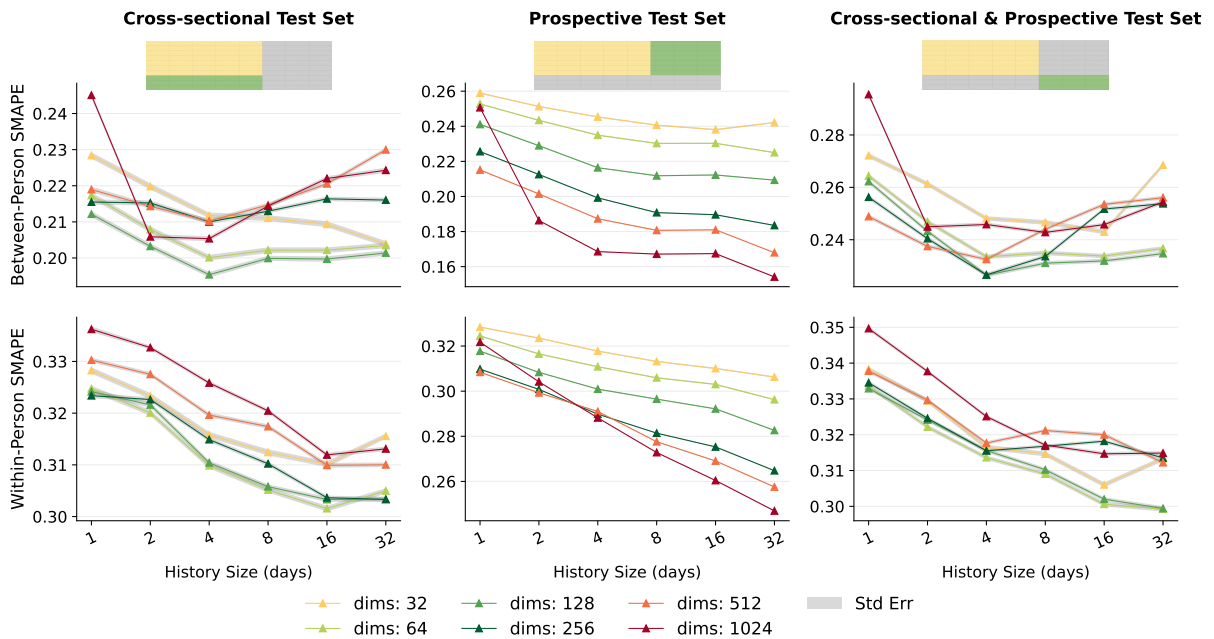


Figure 7: Between- and Within-Person metrics as a function of history length for Bag of Embeddings Model. Predictive Performance improves with modeling the temporal context of linguistic behavior. For Prospective test set, between- and within-person performance improves with longer temporal context and higher dimensional sizes (512 and 1024). For Cross-sectional and Cross-sectional & Prospective evaluation sets, performance improves with temporal context at lower dimensions (hidden size=64-128).

Figure 8: Comparison of Auto Regressive, Bag of Embeddings and Transformer model using Between- (top) and Within-Person (bottom) SMAPE as a function of history length.

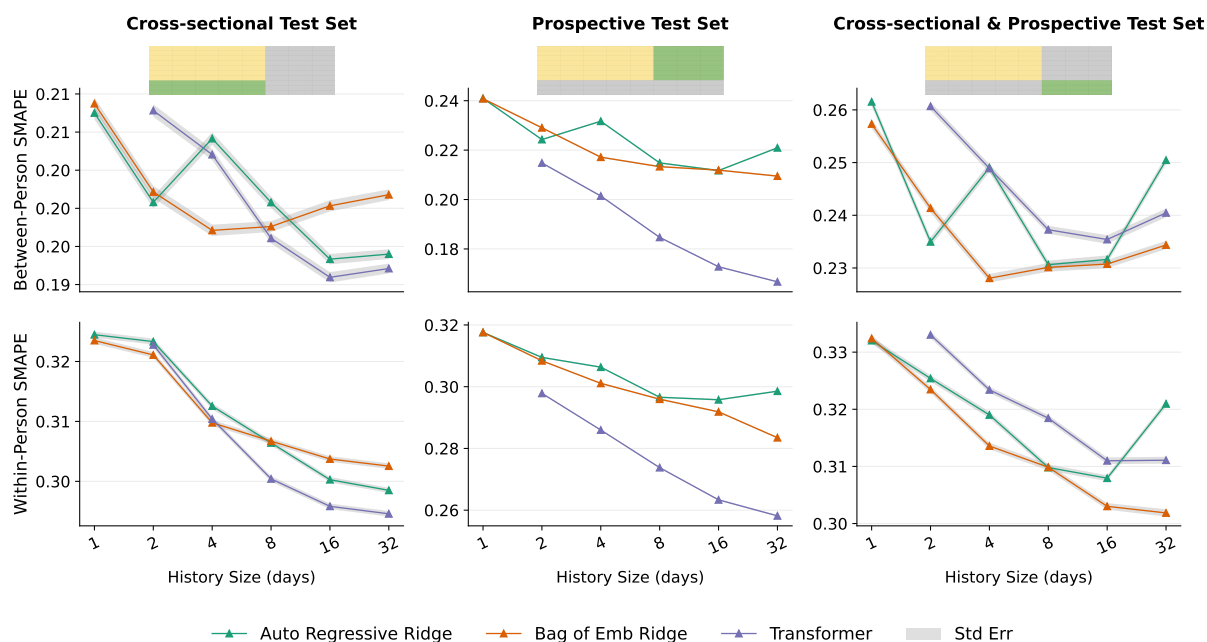


Figure 9: **AR vs. BoE vs. Transformer across history length (128 dims per day).** Between-person (top) and within-person (bottom) SMAPE as a function of history h . Modeling temporal interactions (Transformer) offers marginal improvements over modeling dynamics in cross-sectional splits, but yields the largest gains for prospective generalization. It performs worst for cross-sectional & prospective.