HOW TO MODEL HUMAN ACTIONS DISTRIBUTION WITH EVENT SEQUENCE DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper studies forecasting of the future distribution of events in human action sequences, a task essential in domains like retail, finance, healthcare, and recommendation systems where the precise temporal order is often less critical than the set of outcomes. We challenge the dominant autoregressive paradigm and investigate whether explicitly modeling the future distribution or order-invariant multi-token approaches outperform order-preserving methods. We analyze local order invariance and introduce a KL-based metric to quantify temporal drift. We find that a simple explicit distribution forecasting objective consistently surpasses complex implicit baselines. We further demonstrate that mode collapse of predicted categories is primarily driven by distributional imbalance. This work provides a principled framework for selecting modeling strategies and offers practical guidance for building more accurate and robust forecasting systems.

1 Introduction

In many real-world prediction tasks, the precise temporal ordering of events is irrelevant. Instead, predicting the distribution of outcomes, where only the presence or absence of specific elements matters, is sufficient and often more practical.

For instance, in retail operations, probabilistic demand forecasting enables optimal inventory management and supply chain planning by modeling the full range of possible product demands without requiring sequence order (Nassibi et al., 2023; Larson, 2001). Similarly, in healthcare, clinical diagnosis systems treat disease categories as unordered sets within a single hospital admission. The presence of certain conditions is clinically more significant than the exact order in which they were diagnosed (Johnson et al., 2016; Mullenbach et al., 2018). Recommendation systems further exemplify this principle known as *basket prediction* (Rendle, 2020). Finally, many multi-label problems can naturally be framed as distribution forecasting tasks.

The central focus of this paper is to model the future distribution of human actions over a fixed future horizon. In this work we consider Event Sequences (EvS) (Osin et al., 2025; Udovichenko et al., 2024) - temporal records of human actions which underpin a wide range of decision-making systems across domains including healthcare (Johnson et al., 2016), financial transactions (Udovichenko et al., 2024; Mollaev et al., 2024; Yang & Xu, 2019), e-commerce (Li et al., 2021), recommender systems (Shevchenko et al., 2024; Klenitskiy et al., 2024; Zhelnin et al., 2025), and human action recognition (Surkov et al., 2024). Despite its practical importance and deceptively simple formulation, distribution forecasting for EvS remains significantly understudied.

Inspired by advances in Natural Language Processing (NLP), contemporary approaches to modeling sequential behavior often default to autoregressive generation predicting the next token conditioned on an exact prefix ordering (Karpukhin et al., 2024; Klenitskiy et al., 2024). While Next Token Prediction (NTP) has long dominated sequential modeling, Multi-Token Prediction (MTP) has recently gained traction due to its demonstrated improvements in model quality and generalization, particularly in tasks such as planning, code generation and EvS forecasting (Nagarajan et al., 2025; Bachmann & Nagarajan, 2024; Yu et al., 2025; Karpukhin & Savchenko, 2024).

This raises a practical question: When should we model future event distributions **explicitly**, and when is it worth preserving temporal structure through implicit methods like NTP or Multi-Token Prediction (MTP)?

Two challenges complicate this choice. First, event sequences often exhibit *local order invariance*: within short time windows, the precise ordering of actions (e.g., "buy bread" vs. "buy aspirin") may be arbitrary or uninformative, especially in transactional domains (Klenitskiy et al., 2024; Osin et al., 2025; Udovichenko et al., 2024). Second, as we demonstrate empirically in this work, NTP frequently suffers from *mode collapse* on certain datasets. We investigate three hypotheses for this collapse:

- H1 (Global Order Irrelevance): Full sequence order is uninformative.
- **H2** (**Local Order Irrelevance**): Only coarse temporal structure matters; fine-grained order is uninformative. NTP, by overemphasizing local order, may fail to capture these higher-order distributional patterns.
- H3 (Distributional Imbalance): Highly skewed label distributions can independently cause mode
 collapse especially under maximum likelihood training with NTP, which tends to overweight
 frequent tokens.

To resolve these questions, we introduce a diagnostic-driven framework for distribution forecasting:

- Quantifying Temporal Drift: We introduce a KL-based staticity index, which quantifies distributional drift over time ,providing a dataset-level diagnostic for the relevance of global temporal structure.
- Local Order Invariance: We conduct controlled experiments in which we randomly permute events within sliding windows of varying lengths during to training and measure each dataset's sensitivity to local and global order disruption. Moreover, we provide a formal analysis demonstrating that NTP would fail with locally invariant data structure.
- Mode Collapse Analysis: We quantify distribution characteristics for each dataset, such as the
 exponential decay factor. Although we did not observe any correlation between mode collapse and
 distribution skewness, our results indicate connection of skewness and sensetivity to invariance.
- Explicit vs. Implicit Objectives: We systematically compare four training paradigms: (1) Next Token Prediction (NTP), (2) Multi-Token Prediction (MTP) with ordered output, (3) an order-invariant set prediction approach with post-hoc alignment to targets, and (4) a novel explicit distribution forecasting objective that directly models category probabilities without enforcing order. Our experiments show that order-invariant approaches, particularly the simple explicit method significantly outperform order preserving baselines across most domains. Surprisingly, even on textual data where order is traditionally assumed critical the explicit approach remains superior.

We evaluate all methods across two models and seven public datasets spanning recommendation, retail, banking, and natural language domains. By unifying these findings, our work provides actionable practitioner guidance and paves the way for future research in human behavior modeling.

2 RELATED WORK

Architectures for Event Sequences. Modeling user actions sequentially by conditioning on past behavior has become an essential component of modern recommendation pipelines. These approaches effectively adapt ideas from natural language processing (NLP), particularly attention-based architectures (Kang & McAuley, 2018; Sun et al., 2019; Klenitskiy et al., 2024; Mezentsev et al., 2024). However, it remains unclear whether transformer-based architectures are indeed the most suitable for predicting future user actions. In *EBES* Osin et al. (2025) and in *Seq-NAS* Udovichenko et al. (2024), the authors demonstrate that RNN-based architectures outperform transformer-based models on **EvS** classification tasks. Delving deeper into this issue, Karpukhin & Savchenko (2025) investigate the limitations of transformers and proposes several modifications that enable them to surpass RNNs in classification performance. However, as the same work further reveals, these enhancements do not translate to improved performance in forecasting future tokens. In this work, we focus on RNN- and GPT-based architectures, as they remain the most applicable in this domain.

Multi-Token vs. Single-Token Prediction. Multi-Token Prediction (MTP) has recently gained traction due to its demonstrated improvements in model quality and generalization particularly in tasks such as planning, code generation (Nagarajan et al., 2025; Bachmann & Nagarajan, 2024; Yu

et al., 2025). However, a key challenge lies in the common assumption that predicted tokens are conditionally independent Gloeckle et al. (2024).

Teacherless Learning Bachmann & Nagarajan (2024) offers an intermediate approach between Next-Token Prediction (NTP) and MTP, conceptually opposing teacher forcing. Unlike MTP, Teacherless Learning is grounded in a rigorous mathematical framework. While it does not accelerate inference, it addresses fundamental limitations of traditional NTP. As Nagarajan et al. (2025) note: "Teacherless training and diffusion models comparatively excel in producing diverse and original output."

Although earlier work focused primarily on text generation, Karpukhin & Savchenko (2024) extended these ideas to **EvS**, demonstrating that multi-token generation and diffusion-based approaches indeed outperform the single-token paradigm. In this work, we investigate NTP, a multi-token strategy similar to that proposed in Karpukhin & Savchenko (2024) and propose a new explicit approach for distribution forecasting.

Order Importance in EvS. It has been established that permuting sequences in **EvS** datasets does not degrade performance on classification tasks (Osin et al., 2025; Moskvoretskii et al., 2024), an observation which significantly challenges the assumed sequential nature of this data type. Klenitskiy et al. (2024) investigates whether datasets from the domain of sequential recommender systems genuinely exhibit sequential structure. Specifically, the authors evaluate whether permuting sequences leads to performance degradation in next-token prediction tasks, and find that the extent of degradation varies by dataset, some datasets are more "sequential" than others. In this work, we extend this investigation beyond recommender systems and analyze local permutation invariance, as discussed in **H2** (Section 1).

3 DATASETS

To evaluate the proposed methods and hypotheses, we conduct experiments on a diverse collection of real-world sequential datasets spanning multiple domains—including financial transactions, ecommerce, retail, music streaming, and literary text. A summary of key statistics is provided in Table 1; full descriptions, including preprocessing steps are available in Appendix A.3.

Dataset	ID	Domain	Sequences	Mean len	Target Field	Classes
Multimodal Banking Dataset 2024	MBD	Transactions	1.5M	313	Event type	55
AgeGroup Transactions	AGE	Transactions	30K	888	Small group	203
X5 RetailHero	Retail	Retail	40K	112	Level 2	43
Alphabattle-2.0	AB	Transactions	1M	213	MCC category	28
Complete Works of Shakespeare	ShS	Text	5K	106	Character	65
Megamarket (2024)	MM	E-commerce	2.73M	653	Category ID	9.8K
Zvuk (2024)	Zvuk	Music Streaming	380K	1020	Artist ID	210K
Taobao User Behavior	Taobao	E-commerce	10K	535	Item category	8K

Table 1: Dataset statistics and characteristics.

4 DATASET DIAGNOSTIC

4.1 TEMPORAL ORDER AND MODE COLLAPSE IN EVENT SEQUENCE MODELING

In time series and natural language modeling, precise temporal ordering is crucial. However, in domains like system logs or bank transactions, the *exact micro-temporal order* of events within short windows may be ambiguous or irrelevant—e.g., two unrelated log entries milliseconds apart could plausibly appear in either order without changing system semantics. We illustrate this effect in Appendix 4. This motivates a formal distinction between two types of temporal structure:

• Local invariance: Within a narrow window $W_t = (y_t, \dots, y_{t+H})$, event order is semantically irrelevant—permutations of the same multiset are equally plausible.

• Global structure: Across broader time intervals, dependencies between consecutive windows remain meaningful; e.g., $p(W_2 \mid W_1)$ for $W_1 = (y_0, \dots, y_{t-1})$ and $W_2 = (y_t, \dots, y_{t+H})$ captures genuine temporal progression.

Conventional autoregressive (AR) models are trained to predict the next token y_t given its full history (y_0,\ldots,y_{t-1}) . To accommodate local invariance, one might relax this strict left-to-right dependency by defining a *prediction horizon* $\{y_t,\ldots,y_{t+H}\}$ and training the model to predict *any* event within this window. Under the assumption of uniform uncertainty over the horizon, the training objective becomes:

$$\mathbb{E}_{k \sim \text{Uniform}[0,H]} \left[\log p(x = y_{t+k} \mid y_0, \dots, y_{t-1}) \right] = \frac{1}{H+1} \sum_{m=0}^{H} \log p(x = y_{t+m} \mid y_0, \dots, y_{t-1}).$$
(1)

Critically, standard AR architectures use a *single output distribution* $q_t(\cdot)$ at time t to score all tokens in the horizon. Under local permutation invariance, the optimal q_t that maximizes the above objective is the empirical distribution over the multiset $\{y_t,\ldots,y_{t+H}\}$. Consequently, the model learns a *static predictive distribution* over the entire window: $q_t \approx q_{t+1} \approx \cdots \approx q_{t+H}$. This static distribution becomes problematic at inference time. When generating sequences using deterministic decoding (e.g., argmax or low-temperature sampling), the model outputs:

$$\hat{y}_{t+k} = \arg\max_{x} q_{t+k}(x) \approx \arg\max_{x} q_{t}(x), \quad \forall k \in [0, H].$$

Since q_t is dominated by the most frequent event in the window, the model repeatedly predicts the *empirical mode* of W_t , suppressing rarer—but valid—events. We term this phenomenon *temporal mode collapse*.

We propose that explicitly modeling the distribution of events across entire windows, rather than enforcing pointwise predictions, offers a principled resolution. This allows models to better capture the stochastic nature of real-world event sequences while avoiding degenerate solutions.

4.2 STATICITY INDEX

Before fitting neural models, we quantify how each sequence's event distribution changes over time. Several datasets contain sequences with nearly static behaviour; to verify this, we plot the *Shape* score drift for each dataset.

4.2.1 Per-feature dissimilarity score

Procedure. For each sequence, we fix a window length W and stride s, then slide the window across the timeline. At every position i, we extract the feature distribution P_i within the current window and compare it with the baseline distribution P_0 computed from the first window. To compare them we suggest to leverage the following score:

Let P_0 and P_i denote the empirical distributions in the reference and the *i*-th window, respectively.

Discrete features. For categorical attributes defined on \mathcal{A} we employ the *total variation (TV) distance*, $\mathrm{TV}(P_0, P_i) = \frac{1}{2} \sum_{a \in \mathcal{A}} |P_0(a) - P_i(a)|$. Because lower TV indicates higher similarity, we report its complement $(1 - \mathrm{TV})$, so that higher values consistently reflect better alignment.

Continuous features. For numerical attributes we use the *Kolmogorov–Smirnov* statistic. Let F_0 and F_i be the empirical CDFs corresponding to P_0 and P_i . The KS divergence is $\mathrm{KS}(P_0,P_i) = \sup_{x \in \mathbb{R}} |F_0(x) - F_i(x)|$. Analogously, we report the similarity score $1 - \mathrm{KS}$.

Shape score. For window i we propose to compute each feature's distance using the appropriate formula above and then average across all features: $\operatorname{Shape}(P_0, P_i) = \frac{1}{M} \sum_{j=1}^M d_j (P_0^{(j)}, P_i^{(j)})$, where M is the number of features and d_j is TV when the jth feature is categorical, and KS otherwise. Plotting $i \mapsto \operatorname{shape}(P_0, P_i)$ yields the drift curves used throughout this paper.

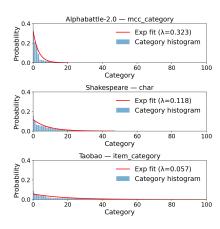


Figure 1: Distribution of categories in datasets. We present normalized number of categories.

Table 2: Dataset statistics: exponential decay parameter (λ) , total number of distinct categories (TCD), and perplexity (PPL) increase after full shuffle.

Dataset	λ	TCD	\mathbf{S}	PPL	
Banking domain					
MBD	0.415	55	0.842	1.02×	
AB	0.305	28	0.725	$1.08 \times$	
Age	0.245	203	0.772	$1.24 \times$	
Retail	0.185	43	0.782	$1.27 \times$	
Text					
ShS	0.118	64	0.803	5.09×	
Recommender Systems					
Taobao	0.016	1.9K	0.650	13.00×	
MM	0.005	9.8K	0.406	$10.87 \times$	
Zvuk	0.003	210K	0.363	$5.05 \times$	

With these definitions we shift the window across the entire sequence and plot trajectory $i \mapsto \operatorname{shape}(P_0, P_i)$, obtaining time-resolved drift curve that summarises how the distribution evolves over the time.

4.2.2 STATICITY IN DATASETS

Across banking datasets (MBD, Retail, Age, AlphaBattle) the majority of user sequences form static clusters with negligible temporal drift (Figure 5, Appendix A.4). In contrast, RecSys data such as ZVUK exhibit more diverse and volatile trajectories (Figure 6, Appendix A.4), while the Shakespeare dataset, despite being textual, resembles banking data with largely flat drift patterns (Figure 7, Appendix A.4). Detailed analyses for individual datasets are provided in the Appendix A.4.

Motivation. These observations motivate a prevent-level *staticity index* that can be computed *before* model training to guide the choice of modeling strategy. Unlike the single–anchor variant (first window vs. all others), we adopt a more robust, multi–anchor formulation.

Staticity index. Fix a window length W and stride s. For each sequence u with per-window distributions $\{P_i^{(u)}\}_{i=1}^{I_u}$, choose anchors \mathcal{R}_u (uniformly at random, R=3). The per-sequence score is the average shape-similarity

$$S^{(u)} = \frac{1}{RI_u} \sum_{r \in \mathcal{R}_u} \sum_{i=1}^{I_u} \text{Shape}(P_r^{(u)}, P_i^{(u)}),$$

and the dataset–level index is Staticity = $\frac{1}{N} \sum_{u=1}^{N} S^{(u)}$.

Thus, the staticity index quantifies the temporal stability of a sequence's multi-feature distribution: higher values (close to 1) reflect stronger staticity (quasi-stationarity), whereas values near zero indicate pronounced drift. Importantly, the conclusions derived from the computed staticity index (Table 2) align with those previously inferred from the qualitative analysis of the plots.

4.3 LOCAL PERMUTATION OF EVENTS

To assess the role of temporal order, we permute events within a symmetric window of some size centered at each position. We evaluate different windows, where $w = \{0, 1, 4, 16, -1\}$ denotes a number of permuted neighbors, w = -1 corresponds to permutation of the full sequence. This design reveals whether models rely on local ordering or global sequence structure.

5 DISTRIBUTION FORECASTING METHODS

We study the task of forecasting a distribution of a sequence over some horizon N given its history. To this end, we consider several training objectives — autoregressive, target-based, matched, and our order-invariant formulation. For all experiments N is fixed as 32.

5.1 AUTOREGRESSIVE LOSS

Let $x_{1:T}$ be a sequence with $x_t \in \{1, \dots, K\}$. The model parameterises conditional next-event probabilities $p_{\theta}(x_{t+1} \mid x_{1:t})$ given the preceding context $x_{1:t}$. The sequence log-likelihood factorises as: $\log p_{\theta}(x_{1:T}) = \sum_{t=1}^{T} \log p_{\theta}(x_{t+1} \mid x_{1:t})$.

5.2 TARGET LOSS

In this setting the model predicts an entire block of L future events in a single forward pass, using a fixed prefix $x_{1:T}$ as context; no teacher forcing is applied inside the horizon. Let $\hat{p}_{T+1},\ldots,\hat{p}_{T+L}$ be the categorical distributions produced for positions T+1 through T+L. The target loss is the sum of negative log-likelihoods for that block: $\mathcal{L}_{\text{target}}^{(L)} = \sum_{i=T+1}^{T+L} -\log \hat{p}_i(x_i \mid x_{1:T})$

Unlike the autoregressive objective, every term is conditioned on *the same* prefix $x_{1:T}$; the model **GRU-Target** therefore learns to produce an entire horizon coherently without receiving the ground-truth events $x_{T+1:T+L-1}$ as intermediate inputs.

5.3 MATCHED LOSS

When the temporal order of future events is weakly informative, forcing the model to predict both the *events* and their *exact positions* needlessly penalises near-correct outputs. The **GRU-Matched** model adapts the matching idea of Karpukhin & Savchenko (2024), aligning each target event with the nearest prediction within a tolerance window of size m, treated as a hyperparameter.

Let a fixed prefix $x_{1:T}$ condition a one-shot block prediction $\hat{p}_{T+1:T+L}$; let $x_{T+1:T+L}$ be the corresponding ground truth. With a permutation σ constrained by $|\sigma(i)-i| \leq m$, the matched loss is $\mathcal{L}_{\text{match}}^{(m)} = \min_{\substack{\sigma \in \mathcal{A} \\ |\sigma(i)-i| \leq m}} \sum_{i=T+1}^{T+L} -\log \hat{p}_{\sigma(i)}\big(x_i \mid x_{1:T}\big).$

At m=0 it reduces to plain block cross-entropy; as m grows, the objective becomes progressively order-invariant. The minimisation is solved with the Hungarian algorithm on the cost matrix $\ell_{ij} = -\log \hat{p}_i(x_i \mid x_{1:T})$.

5.4 ORDER-INVARIANT DISTRIBUTION PARAMETERIZATION

When the order of future events is not informative, it is sufficient to model only the *event type distribution* rather than their precise temporal arrangement. We therefore introduce the **GRU-Dist** model, which represents each sequence as a *bag of events* and is trained to match the empirical distribution.

Let $H_t = \{x_1, \dots, x_t\}$ be the multiset of events observed so far. A neural encoder f_θ maps H_t to logits, which are converted to probabilities $\pi_t = \operatorname{softmax} \big(f_\theta(H_t)\big) \; ; \in \; \Delta^{K-1}$, where Δ^{K-1} is the probability simplex in \mathbb{R}^K . For a sequence of length L we form its empirical distribution $\hat{p}_k = \frac{1}{L} \sum_{t=1}^L \mathbf{1}\{x_t = k\}$, and minimize $\ell(\theta) = \operatorname{D}_{\mathrm{KL}} \big(\hat{p} \, \| \, \boldsymbol{\pi}(\theta)\big)$.

Unlike autoregressive objectives that require $L \times K$ logits per sequence, our order-invariant head outputs only a single K-dimensional vector. This reduces both computational and memory costs by a factor of L, while remaining well suited for datasets where event order carries little information.

5.5 MULTI-TOKEN PREDICTION VIA SAMPLING

Autoregressive decoding with greedy argmax often collapses to the modal category. A simple remedy is to *sample* from the predictive categorical distribution instead of always taking the maximum,

which reduces mode collapse and improves order-invariant metrics. For autoregressive and block-prediction models this sampling is straighforward, as logits at each step define the distribution, in our order-invariant method the distribution itself is parameterized directly, making sampling the natural decoding mechanism. We did not analyze more sophisticated sampling approaches such as beam search and our preliminary experiments with temperature sampling did not provide stable improveent across datasets, so we do not use them. **Sampling in the order-invariant model:** Given a predicted categorical distribution $\pi = (\pi_1, \dots, \pi_K)$ and a target length L, we compute expected counts $n_k = L \cdot \pi_k, \sum_{k=1}^K n_k = L$. The category counts are computed using Hamilton's method from apportionment theory Balinski & Young (2010), which distributes L discrete slots among categories in proportion to their predicted probabilities π_k and guarantees that the total count equals L.

6 EVALUATION

For each configuration $Dataset \times Method \times LocalShuffle$ we perform an extensive hyperparameter optimization of 100 trails, technical details are given in Appendix A.1.

6.1 Baselines

We consider four simple baselines. (1) Ground Truth uses the original sequences as a sanity check and reference point for metrics such as Cardinality. Repeat extends a sequence by copying its most recent observations into the forecast horizon of the lenght N. Mode outputs the users most frequent category for all N, illustrating the tendency of autoregressive models to collapse into trivial mode repetition—a behavior that may be overestimated by order-dependent metrics (e.g., Accuracy, Levenshtein distance). Finally, **HistSampler** generates sequences by sampling from the empirical histogram of past users sequence, thereby preserving marginal category frequencies while discarding temporal dependencies.

6.2 METRICS

Many classical sequence metrics (e.g., Accuracy, Levenshtein distance, F1-score) are defined with respect to a fixed token order and therefore penalize any permutation of events, even when such reordering is irrelevant for the problem at hand. To overcome this limitation, we introduce an order-invariant *Matched-F1* score, which treats sequences as *bags of events*.

To avoid order dependence we redefine true-positive, false-positive and false-negative terms. Let g_k and \hat{g}_k denote the ground-truth and predicted multiplicities of class k in the window. We set

$$(TP_k, FP_k, FN_k) = (\min(g_k, \hat{g}_k), \max(0, \hat{g}_k - g_k), \max(0, g_k - \hat{g}_k)).$$

Based on this definitions, we compute our *Matched-F1* with **micro-** and **macro-**averaging, analogous to the conventional F1 formulation. Detailed definition of this metric placed in Appendix A.6.1

To assess diversity, we use **Cardinality** (see Appendix A.6.2), which measures the number of distinct categories generated by the model. Low values signal mode collapse, while values close to the ground-truth indicate faithful event variety. For completeness, we also report **Levenshtein distance**, an order-sensitive metric that, although less relevant to our setting, provides a complementary reference for order preserving methods.

7 RESULTS

Dataset-level statistics. The staticity index serve as useful diagnostics for anticipating whether sequence order is relevant. Results are presented in Table 2. In banking datasets, a single modal category dominates—accounting for more than 50% of all events—leading to high values of both λ and the staticity index. This dominance is also associated with a pronounced performance drop under local permutations, suggesting limited reliance on sequential order.

Table 3: Next N tokens forecasting. *Matched-F1 (micro)* for all datasets and methods including baselines. † denotes sampled version of method.

Method	MBD	Age	AB	Retail	ShS	Taobao	MM	Zvuk
GT	1.000	1.000	1.000	1.000	1.000	0.926	1.000	1.000
Mode	0.520	0.331	0.380	0.219	0.158	0.117	0.156	0.113
Repeat	0.830	0.680	0.700	0.661	0.587	0.257	0.318	0.274
HistSampler	0.804	0.632	0.680	0.640	0.533	0.197	0.244	0.226
GRU	0.528	0.477	0.375	0.207	0.596	0.222	0.250	0.148
GRU^\dagger	0.771	0.628	0.641	0.609	0.596	0.146	0.171	0.126
GPT	0.524	0.476	0.373	0.212	0.594	0.223	0.250	0,192
GPT^\dagger	0.776	0.627	0.629	0.611	0.603	0.151	0.188	0,174
GRU-Target	0.541	0.370	0.403	0.398	0.299	0.196	0.267	0.143
GRU-Target [†]	0.808	0.633	0.670	0.641	0.572	0.154	0.201	0.140
GRU-Matched	0.847	0.704	0.676	0.708	0.688	0.203	0.272	0.202
GRU-Matched [†]	0.827	0.653	0.647	0.667	0.634	0.155	0.203	0.134
GRU-Dist	0.856	0.725	0.736	0.719	0.705	0.178	0.247	0.239

Local permutation experiments further corroborate these findings; results are shown in Figure 3. Shakespeare and Zvuk exhibit sharp performance degradation when sequences are shuffled, indicating strong local sequential structure. In contrast, most banking datasets show little to no degradation, reflecting the irrelevance of event order. This trend is especially evident in Figure 2, which illustrates minimal perplexity degradation under shuffling for these datasets.

Model performance. The order-invariant model *GRU-Dist* achieves the best overall performance on most datasets, with the notable exceptions of Taobao and Megamarket. Further analysis reveals that these two datasets contain long-horizon repetitive patterns of identical events. This characteristic aligns with the observation that the **repeat** baseline performs best on them, as it effectively exploits such redundancy.

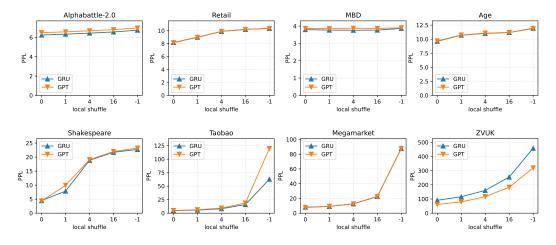


Figure 2: Next N tokens forecasting. Perplexity results.

8 Conclusion

Our study demonstrates that model performance in event-sequence forecasting is tightly linked to dataset characteristics. *GRU-Target* performs best when temporal order is largely irrelevant, while *GRU-Matched* and *GRU-Dist* consistently outperform other approaches; in particular, *GRU-Dist* is considered as more appropriate baseline for banking tasks, where the presence of events is more

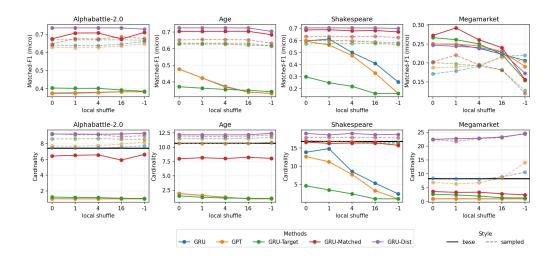


Figure 3: Effect of Local Event Shuffling on Model Performance. We report Matched-F1 score and Carnality for four datasets. Results for other datasets and metrics can be found in Appendix A.8

informative than their order. The drop in GRU-Dist performance on Megamarket and Taobao is likely due to dataset-specific characteristics: sequences often contain long run of identical tokens, which is difficult to reproduce when sampling from a categorical distribution. GPT-based models, by contrast, are more effective on datasets sensitive to local permutations, such as text. Therefore, in scenarios where temporal order is essential, Next-Token Prediction (NTP) and Multi-Token Prediction remain the preferable approaches.

Cardinality also proves to be a useful diagnostic of mode collapse: in datasets like Shakespeare, shuffling removes structural cues and autoregressive models degenerate to the modal category. More broadly, when no meaningful local ordering exists, models tend to collapse to the mode (Figure 3).

Taken together, these results highlight the value of simple dataset-level diagnostics (exponential decay parameter λ , staticity, cardinality) for anticipating model behavior, and demonstrate the advantages of order-invariant objectives in domains such as retail and banking, where event presence matters more than sequence order.

Indeed, it is worth noting that the proposed *GRU-Dist* method can be extended from single-category forecasting to multi-feature prediction through cascade modeling.

Acknowledgment on LLM assisted writing: This paper used open access Qwen3-Max, in some parts of the paper, for proofreading and text rephrasing in accordance with formal style.

REFERENCES

Gregor Bachmann and Vaishnavh Nagarajan. The pitfalls of next-token prediction. *arXiv preprint arXiv:2403.06963*, 2024.

Michel L Balinski and H Peyton Young. Fair representation: meeting the ideal of one man, one vote. Rowman & Littlefield. 2010.

Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*, 2024.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.

Wang-Cheng Kang and Julian McAuley. Self-attentive sequential recommendation. In 2018 IEEE international conference on data mining (ICDM), pp. 197–206. IEEE, IEEE Computer Society, 2018.

- Ivan Karpukhin and Andrey Savchenko. Detpp: Leveraging object detection for robust long-horizon event prediction. arXiv preprint arXiv:2408.13131, 2024.
- Ivan Karpukhin and Andrey Savchenko. Ht-transformer: Event sequences classification by accumulating prefix information with history tokens. *arXiv* preprint arXiv:2508.01474, 2025.
 - Ivan Karpukhin, Foma Shipilov, and Andrey Savchenko. Hotpp benchmark: Are we good at the long horizon events forecasting? *arXiv preprint arXiv:2406.14341*, 2024.
 - Anton Klenitskiy, Anna Volodkevich, Anton Pembek, and Alexey Vasilev. Does it look sequential? an analysis of datasets for evaluation of sequential recommendations. In *Proceedings of the 18th ACM Conference on Recommender Systems*, pp. 1067–1072, 2024.
 - Paul D Larson. Designing and managing the supply chain: concepts, strategies, and case studies. *Journal of Business Logistics*, 22(1):259, 2001.
 - LinShu Li, Jianbo Hong, Sitao Min, and Yunfan Xue. A novel ctr prediction model based on deepfm for taobao data. In 2021 IEEE International Conference on Artificial Intelligence and Industrial Design (AIID), pp. 184–187. IEEE, 2021.
 - Gleb Mezentsev, Danil Gusak, Ivan Oseledets, and Evgeny Frolov. Scalable cross-entropy loss for sequential recommendations with large item catalogs. In *Proceedings of the 18th ACM Conference on Recommender Systems*, pp. 475–485, 2024.
 - Dzhambulat Mollaev, Alexander Kostin, Maria Postnova, Ivan Karpukhin, Ivan Kireev, Gleb Gusev, and Andrey Savchenko. Multimodal banking dataset: Understanding client needs through event sequences. *arXiv preprint arXiv:2409.17587*, 2024.
 - Viktor Moskvoretskii, Dmitry Osin, Egor Shvetsov, Igor Udovichenko, Maxim Zhelnin, Andrey Dukhovny, Anna Zhimerikina, Albert Efimov, and Evgeny Burnaev. Self-supervised learning in event sequences: A comparative study and hybrid approach of generative modeling and contrastive learning. *CoRR*, 2024.
 - Jeff Mullenbach, Ross Swanson, James Wallis, Hila Bekerman, Michael Chiang, and Jim Glass. Explainable prediction of medical codes from clinical text. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 1, pp. 1101–1111, 2018.
 - Vaishnavh Nagarajan, Chen Henry Wu, Charles Ding, and Aditi Raghunathan. Roll the dice & look before you leap: Going beyond the creative limits of next-token prediction. *arXiv* preprint *arXiv*:2504.15266, 2025.
 - Nouran Nassibi, Heba Fasihuddin, and Lobna Hsairi. Demand forecasting models for food industry by utilizing machine learning approaches. *International Journal of Advanced Computer Science and Applications*, 14(3):892–898, 2023.
 - Dmitry Osin, Igor Udovichenko, Egor Shvetsov, Viktor Moskvoretskii, and Evgeny Burnaev. Ebes: Easy benchmarking for event sequences. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 5730–5741, 2025.
 - Steffen Rendle. Recommender systems. *Foundations and Trends in Information Retrieval*, 14(1-2): 1–223, 2020.
 - Valeriy Shevchenko, Nikita Belousov, Alexey Vasilev, Vladimir Zholobov, Artyom Sosedka, Natalia Semenova, Anna Volodkevich, Andrey Savchenko, and Alexey Zaytsev. From variability to stability: Advancing recsys benchmarking practices. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 5701–5712, 2024.
 - Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pp. 1441–1450, 2019.

Egor Surkov, Oleg Seredin, and Andrei Kopylov. Human action recognition based on the skeletal pairwise dissimilarity. *Computer Optics*, 2024.

Igor Udovichenko, Egor Shvetsov, Denis Divitsky, Dmitry Osin, Ilya Trofimov, Ivan Sukharev, Anatoliy Glushenko, Dmitry Berestnev, and Evgeny Burnaev. Seqnas: Neural architecture search for event sequence classification. *IEEE Access*, 12:3898–3909, 2024.

Kunlin Yang and Wei Xu. Fraudmemory: Explainable memory-enhanced sequential neural networks for financial fraud detection. 2019.

Tianhao Yu, Xianghong Zhou, and Xinrong Deng. Autoregressive models for session-based recommendations using set expansion. *PeerJ Computer Science*, 11:e2734, 2025.

Maxim Zhelnin, Dmitry Redko, Volkov Daniil, Anna Volodkevich, Petr Sokerin, Valeriy Shevchenko, Egor Shvetsov, Alexey Vasilev, Darya Denisova, Ruslan Izmailov, et al. Faster and memory-efficient training of sequential recommendation models for large catalogs. *arXiv* preprint arXiv:2509.09682, 2025.

A APPENDIX

A.1 HPO DETAILS

For hyperparameter optimization (HPO), we use Optuna (?) with the Tree-structured Parzen Estimator (TPE) sampler. For each model-dataset pair, we allocate an HPO budget of 100 training runs, capping the total computational cost at 18 NVIDIA A100 GPU-days. We reserve 15% of the training set as a validation subset for early stopping and hyperparameter selection. The best-performing hyperparameters are then used to train the final model for evaluation and all subsequent study experiments.

A.2 LOCAL GLOBAL TEMPORAL INVARIANCE

In Figure 4 we illustrate local / global invariance.

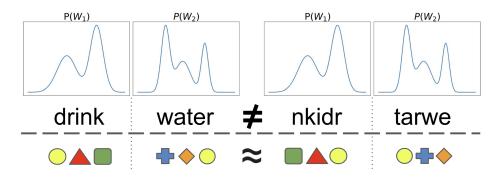


Figure 4: Example how order importance differs in different types of data. Even though in both cases horizon distribution doesnt change, event sequence still make sence after permut inside intervals.

A.3 DATASETS DESCRIPTION AND PREPROCESSING

MBD ¹ is a multimodal banking dataset introduced in ?. The dataset contains an industrial-scale number of sequences, with data from more than 1.5 million clients in 2 year period. Each client corresponds to a sequence of events. This multi-modal dataset includes card transactions, geoposition events, and embeddings of dialogs with technical support. For our analysis, we use only card transactions. We use a temporal train—test split: transactions from the first year form the training set, and those from the second year form the test set.

https://huggingface.co/datasets/ai-lab/MBD

Age dataset² consists of 44M anonymized credit card transactions representing 50K individuals. The target is to predict the age group of a cardholder that made the transactions. Each transaction includes the date, type, and amount being charged. The dataset was first introduced in scientific literature in work ?. We perform a user-based split: 80% of sequences are assigned to the training set, and the remaining 20% of sequences are held out for testing.

Retail dataset³ comprises 45.8M retail purchases from 400K clients, with the aim of predicting a client's age group based on their purchase history. Each purchase record includes details such as time, item category, the cose, and loyalty program points received. The age group information is available for all clients, and the distribution of these groups is balanced across the dataset. The dataset was first introduced in scientific literature in work? We perform a user-based split: 80% of sequences are assigned to the training set, and the remaining 20% of sequences are held out for testing.

Alphabattle-2.0 datase ⁴ The AlfaBattle 2.0 dataset contains bank customers' transaction records over two years, with the goal of predicting loan default based on behavioral history. Each record includes 18 features (3 numeric, 15 categorical) per transaction. We use the official test split provided by the dataset creators.

Shakespeare Dataset consists of character-level text extracted from Shakespeare's works, preprocessed into individual speech segments. Each speech is tokenized using a vocabulary of unique characters mapped to integer codes based on frequency. The final dataset is split into train and test sets (80/20). The dataset is designed for character-level language modeling and was selected due to it obvious temporal importance.

Zvuk dataset⁵ is introduced in 2024 and contains 244.7M music listening events grouped into 12.6M sessions from 382K users, recorded during the same five-month period (January–May 2023). In total, it spans 1.5M unique tracks. Each record includes a user ID, session ID, track ID, timestamp, and play duration (considering only plays covering at least 30% of track length). The dataset is tailored to music consumption, excluding podcasts and audiobooks, and enables evaluation of recommendation models in domains with stronger sequential dynamics. We use a temporal train–test split: transactions from the first two months form the training set, and other two month form the test set.

MegaMarket dataset⁶ is introduced in 2024 and comprises 196.6M user interactions collected over a five-month period (January–May 2023). It covers 2.7M users, 3.56M items, and 10,001 product categories, with events including views, favorites, cart additions, and purchases. Each record contains a user ID, item ID, event type, category ID, timestamp, and normalized price. The dataset represents large-scale e-commerce behavior and is intended for sequential recommendation tasks. This dataset follows the same temporal train/test split as Zvuk.

Taobao ⁷ The dataset comprises user behaviors from Taobao, including clicks, purchases, adding items to the shopping cart, and favoriting items. These events were collected between November 18 and December 15. The training set encompasses data from November 18 to December 1, while the test set includes clicks from December 2 to December 15.

A.4 STATICITY INDEX PLOTS FOR KEY DATASETS

⁷https://tianchi.aliyun.com/dataset/46

For each dataset, we compute drift trajectories for all sequence and cluster them into a small number of groups with internally consistent dynamics (Figure 5–7). Across banking datasets (MBD, Retail,

```
2https://ods.ai/competitions/sberbank-sirius-lesson
3https://ods.ai/competitions/x5-retailhero-uplift-modeling
4https://www.kaggle.com/datasets/mrmorj/alfabattle-20
5https://www.kaggle.com/datasets/alexxl/zvuk-dataset
6https://www.kaggle.com/datasets/alexxl/megamarket?select=megamarket.parquet
```

Age, Alphabattle) the dominant clusters are static, as exemplified for **MBD** (Figure 5c), these clusters exhibit negligible temporal drift. For such sequences, learning the user's category distribution suffices to forecast the next block of events. Trajectories with pronounced drift are rare. In MBD specifically, such sequences are observed in fewer than 6% of users (Figure 5b).

In contrast to banking datasets, recommender–system data exhibit much greater variability. In **ZVUK** (Figure 6), two characteristic regimes dominate: one cluster shows a sharp initial drop from the baseline followed by persistent high-variance fluctuations, while another appears quasi-static yet remains noisy around its trend. Such patterns reflect the broader nature of recommender logs: users interact with a large and diverse sets of items, and their behavior shifts more frequently than in retail domains where event types are limited and highly regular. And as a consequence, their later-window distributions are more clearly separated from the first-window distribution.

The outlier in this collection is the **Shakespeare** text dataset (Figure 7). Although it is non-transactional, its dynamics resemble banking data more than recommender logs: drift trajectories are mostly flat and volatility remains low. At the same time, weak periodic or gradual shifts are observable, indicating that the sequences are not fully static but display a modest degree of temporal variation.

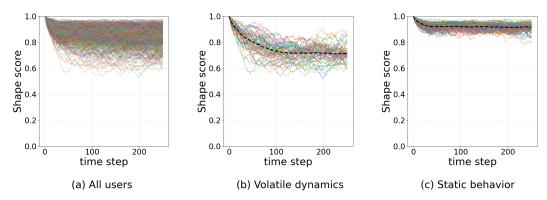


Figure 5: Shape score drift for MBD dataset

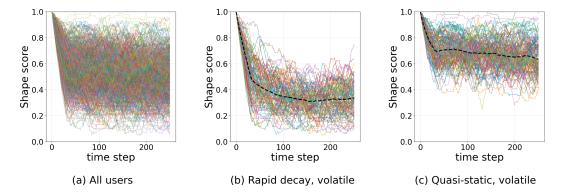


Figure 6: Shape score drift for ZVUK dataset

A.5 FEATURES IMPACT IN CATEGORY FORECASTING QUALITY

We investigated whether predicting a target feature benefits more from incorporating the full feature vector or from relying exclusively on its own historical values.

On the MBD dataset, experiments in the *All-to-One* and *One-to-One* modes reveal that the autoregressive model's performance degrades when exposed to complete with the complete feature vector. The additional inputs act as noise, impeding the model's ability to reproduce the mode of the target distribution. In the *One-to-One* mode—where the model sees only the history of the target

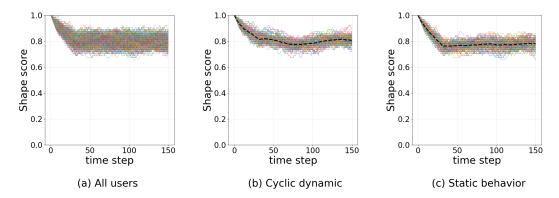


Figure 7: Shape score drift for Shakespeare dataset

feature—it easily learns the mode and reports a formal increase in accuracy; however, this gain is illusory, as the generated sequences become overly uniform and lack realism 4.

Table 4: Effect of training with all tokens vs. event type only (*Matched-F1 micro*).

Dataset	Change (%)
MBD	+2.85
AGE	-24.94
MM	+13.66

By contrast, on datasets with a strong sequential structure, such as *Megamarket*, the opposite pattern emerges. The autoregressive mechanism leverages ordering information and, when augmented with additional features, predicts beyond mere modal values, resulting in a significant improvement in performance metrics.

A.6 METRICS

A.6.1 MATCHED-F1 MICRO

Precision and recall.

$$\mathrm{Prec}_k = \frac{\mathrm{TP}_k}{\mathrm{TP}_k + \mathrm{FP}_k}, \qquad \mathrm{Rec}_k = \frac{\mathrm{TP}_k}{\mathrm{TP}_k + \mathrm{FN}_k}.$$

Macro averaging.

$$F1_{\text{macro}} = \frac{1}{K} \sum_{k=1}^{K} \frac{2 \operatorname{Prec}_{k} \operatorname{Rec}_{k}}{\operatorname{Prec}_{k} + \operatorname{Rec}_{k}}.$$

Each class contributes equally; the score is sensitive to rare categories.

Micro averaging. Aggregating counts over classes,

$$TP = \sum_{k} TP_{k}, \qquad FP = \sum_{k} FP_{k}, \qquad FN = \sum_{k} FN_{k}, \qquad (2)$$

$$F1_{\text{micro}} = \frac{2 \text{ TP}}{2 \text{ TP} + \text{FP} + \text{FN}}.$$
 (3)

This variant weights categories by frequency and reflects overall throughput.

A.6.2 CARDINALITY METRIC.

Let $G_i = (x_{t+1}^{(i)}, \dots, x_{t+L}^{(i)})$ denote the L-step segment generated for sequence i and $\mathcal{C}(G_i) = \{x \in G_i\}$ the set of *distinct* categories appearing in that segment. We define the per-sequence cardinality as

$$C_i = |\mathcal{C}(G_i)|.$$

The dataset-level score is the average

Cardinality =
$$\frac{1}{N} \sum_{i=1}^{N} C_i$$
,

where N is the number of sequences under evaluation. An *overall* variant first concatenates all generated segments, $\tilde{G} = \bigcup_i G_i$, and reports $C_{\text{overall}} = |\mathcal{C}(\tilde{G})|$.

Purpose. Cardinality captures the *category diversity* produced by a model: low values signal mode collapse, whereas values close to the ground-truth cardinality indicate faithful reproduction of event variety. We compute the metric for both generated $(C_{\rm gen})$ and reference $(C_{\rm orig})$ sequences, allowing direct comparison of a model's diversity against empirical data.

A.7 NEURAL BACKBONE ARCHITECTURES

We evaluate two neural backbone architectures for sequence modeling:

- **GRU:** A gated recurrent unit (GRU) network excels at capturing local dependencies and stationary patterns in short to moderately long time series.
- **GPT:** A self-attention—based model capable of modeling long-range dependencies, crucial for sequences with complex contextual interactions and implicit event relationships.

A.8 ADDITIONAL RESULTS

For completeness, we report all evaluation metrics across datasets. Levenshtein distance is included as an order-sensitive measure to quantify degradation under local shuffling (Figure 8), while the effect of shuffling on category diversity is illustrated by cardinality (Figure 9). The main text focuses on the order-invariant *Matched-F1 (micro)* (Figure 10), which we adopt as the primary evaluation metric throughout the study.

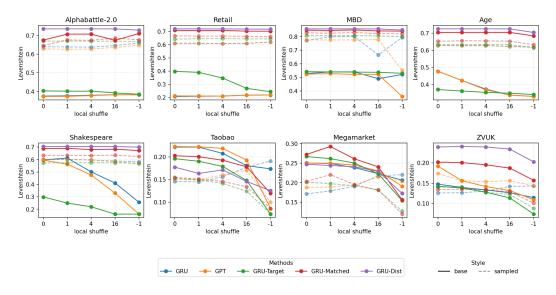


Figure 8: Levenshtein score on all datasets.

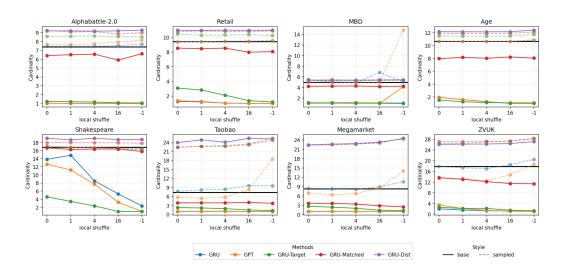


Figure 9: Effect of local shuffle on cardinality.

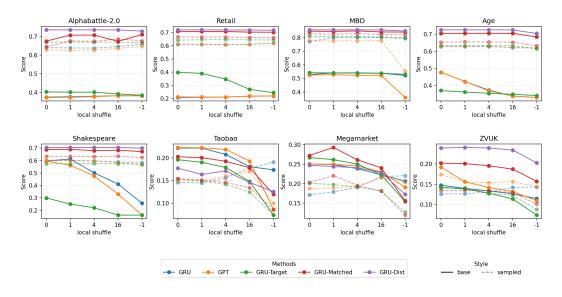


Figure 10: Next N tokens forecasting. Matched-F1 (micro) results.