

MotionLM: Multi-Agent Motion Forecasting as Language Modeling

Ari Seff Brian Cera Dian Chen* Mason Ng Aurick Zhou Nigamaa Nayakanti
 Khaled S. Refaat Rami Al-Rfou Benjamin Sapp

Waymo

Abstract

Reliable forecasting of the future behavior of road agents is a critical component to safe planning in autonomous vehicles. Here, we represent continuous trajectories as sequences of discrete motion tokens and cast multi-agent motion prediction as a language modeling task over this domain. Our model, MotionLM, provides several advantages: First, it does not require anchors or explicit latent variable optimization to learn multimodal distributions. Instead, we leverage a single standard language modeling objective, maximizing the average log probability over sequence tokens. Second, our approach bypasses post-hoc interaction heuristics where individual agent trajectory generation is conducted prior to interactive scoring. Instead, MotionLM produces joint distributions over interactive agent futures in a single autoregressive decoding process. In addition, the model’s sequential factorization enables temporally causal rollouts. The proposed approach establishes new state-of-the-art performance for multi-agent motion prediction on the Waymo Open Motion Dataset, ranking 1st on the interactive challenge leaderboard.

1. Introduction

Modern sequence models often employ a next-token prediction objective that incorporates minimal domain-specific assumptions. For example, autoregressive language models [3, 10] are pre-trained to maximize the probability of the next observed subword conditioned on the previous text; there is no predefined notion of parsing or syntax built in. This approach has found success in continuous domains as well, such as audio [2] and image generation [49]. Leveraging the flexibility of arbitrary categorical distributions, the above works represent continuous data with a set of discrete

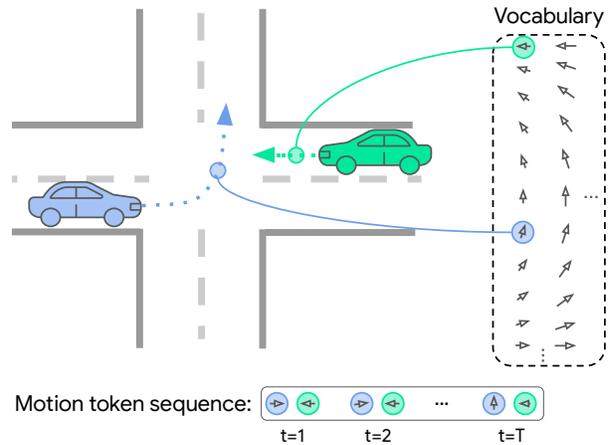


Figure 1. Our model autoregressively generates sequences of discrete motion tokens for a set of agents to produce consistent interactive trajectory forecasts.

tokens, reminiscent of language model vocabularies.

In driving scenarios, road users may be likened to participants in a constant dialogue, continuously exchanging a dynamic series of actions and reactions mirroring the fluidity of communication. Navigating this rich web of interactions requires the ability to anticipate the likely maneuvers and responses of the involved actors. Just as today’s language models can capture sophisticated distributions over conversations, can we leverage similar sequence models to forecast the behavior of road agents?

A common simplification to modeling the full future world state has been to decompose the joint distribution of agent behavior into independent per-agent marginal distributions. Although there has been much progress on this task [8, 47, 12, 25, 31, 5, 6, 21], marginal predictions are insufficient as inputs to a planning system; they do not represent the future dependencies between the actions of different agents, leading to inconsistent scene-level forecasting.

Of the existing joint prediction approaches, some apply

*Work done during an internship at Waymo.
 Contact: {aseff, bensapp}@waymo.com

a separation between marginal trajectory generation and interactive scoring [40, 42, 29]. For example, Luo et al. [29] initially produce a small set of marginal trajectories for each agent independently, before assigning a learned potential to each inter-agent trajectory pair through a belief propagation algorithm. Sun et al. [42] use a manual heuristic to tag agents as either influencers or reactors, and then pairs marginal and conditional predictions to form joint predictions.

We also note that because these approaches do not explicitly model temporal dependencies within trajectories, their conditional forecasts may be more susceptible to spurious correlations, leading to less realistic reaction predictions. For example, these models can capture the *correlation* between a lead agent decelerating and a trailing agent decelerating, but may fail to infer which one is likely causing the other to slow down. In contrast, previous joint models employing an autoregressive factorization, e.g., [36, 43, 39], do respect future temporal dependencies. These models have generally relied on explicit latent variables for diversity, optimized via either an evidence lower bound or normalizing flow.

In this work, we combine trajectory generation and interaction modeling in a single, temporally causal, decoding process over discrete motion tokens (Fig. 1), leveraging a simple training objective inspired by autoregressive language models. Our model, MotionLM, is trained to directly maximize the log probability of these token sequences among interacting agents. At inference time, joint trajectories are produced step-by-step, where interacting agents sample tokens simultaneously, attend to one another, and repeat. In contrast to previous approaches which manually enforce trajectory multimodality during training, our model is entirely latent variable and anchor-free, with multimodality emerging solely as a characteristic of sampling. MotionLM may be applied to several downstream behavior prediction tasks, including marginal, joint, and conditional predictions.

This work makes the following contributions:

1. We cast multi-agent motion forecasting as a language modeling task, introducing a temporally causal decoder over discrete motion tokens trained with a causal language modeling loss.
2. We pair sampling from our model with a simple rollout aggregation scheme that facilitates weighted mode identification for joint trajectories, establishing new state-of-the-art performance on the Waymo Open Motion Dataset interaction prediction challenge (6% improvement in the ranking joint mAP metric).
3. We perform extensive ablations of our approach as well as analysis of its temporally causal conditional

predictions, which are largely unsupported by current joint forecasting models.

2. Related work

Marginal trajectory prediction. Behavior predictors are often evaluated on their predictions for individual agents, e.g., in recent motion forecasting benchmarks [14, 9, 4, 51, 37]. Previous methods process the rasterized scene with CNNs [8, 5, 12, 17]; the more recent works represent scenes with points and polygraphs and process them with GNNs [6, 25, 47, 22] or transformers [31, 40, 20]. To handle the multimodality of future trajectories, some models manually enforce diversity via predefined anchors [8, 5] or intention points [40, 52, 28]. Other works learn diverse modes with latent variable modeling, e.g., [24].

While these works produce multimodal future trajectories of individual agents, they only capture the marginal distributions of the possible agent futures and do not model the interactions among agents.

Interactive trajectory prediction. Interactive behavior predictors model the joint distribution of agents’ futures. This task has been far less studied than marginal motion prediction. For example, the Waymo Open Motion Dataset (WOMD) [14] challenge leaderboard currently has 71 published entries for marginal prediction compared to only 14 for interaction prediction.

Ngiam et al. [32] models the distribution of future trajectories with a transformer-based mixture model outputting joint modes. To avoid the exponential blow-up from a full joint model, Luo et al. [29] models pairwise joint distributions. Tolstaya et al. [44], Song et al. [41], Sun et al. [42] consider conditional predictions by exposing the future trajectory of one agent when predicting for another agent. Shi et al. [40] derives joint probabilities by simply multiplying marginal trajectory probabilities, essentially treating agents as independent, which may limit accuracy. Cui et al. [11], Casas et al. [7], Girgis et al. [15] reduce the full-fledged joint distribution using global latent variables. Unlike our autoregressive factorization, the above models typically follow “one-shot” (parallel across time) factorizations and do not explicitly model temporally causal interactions.

Autoregressive trajectory prediction. Autoregressive behavior predictors generate trajectories at intervals to produce scene-consistent multi-agent trajectories. Rhinehart et al. [36], Tang and Salakhutdinov [43], Amirloo et al. [1], Salzmann et al. [39], Yuan et al. [50] predict multi-agent future trajectories using latent variable models. Lu et al. [27] explores autoregressively outputting keyframes via mixtures of Gaussians prior to filling in the remaining states. In [18], an adversarial objective is combined with

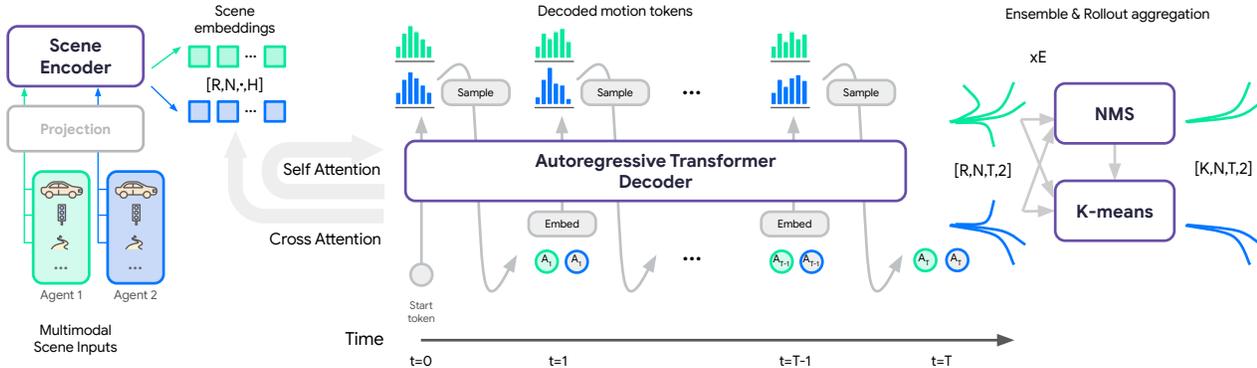


Figure 2. MotionLM architecture. We first encode heterogeneous scene features relative to each modeled agent (left) as scene embeddings of shape R, N, \cdot, H . Here, R refers to the number of rollouts, N refers to the number of (jointly modeled) agents, and H is the dimensionality of each embedding. We repeat the embeddings R times in the batch dimension for parallel sampling during inference. Next, a trajectory decoder autoregressively rolls out T discrete motion tokens for multiple agents in a temporally causal manner (center). Finally, representative modes of the rollouts may be recovered via a simple aggregation utilizing k-means clustering initialized with non-maximum suppression (right).

parallel beam search to learn multi-agent rollouts. Unlike most autoregressive trajectory predictors, our method does not rely on latent variables or beam search and generates multimodal joint trajectories by directly sampling from a learned distribution of discrete motion token sequences.

Discrete sequence modeling in continuous domains.

When generating sequences in continuous domains, one effective approach is to discretize the output space and predict categorical distributions at each step.

For example, in image generation, van den Oord et al. [45] sequentially predict the uniformly discretized pixel values for each channel and found this to perform better than outputting continuous values directly. Multiple works on generating images from text such as [35] and [49] use a two-stage process with a learned tokenizer to map images to discrete tokens and an autoregressive model to predict the discrete tokens given the text prompt. For audio generation, WaveNet [46] applies a μ -law transformation before discretizing. Borsos et al. [2] learn a hierarchical tokenizer/detokenizer, with the main transformer sequence model operating on the intermediate discrete tokens. When generating polygonal meshes, Nash et al. [30] uniformly quantize the coordinates of each vertex. In MotionLM, we employ a simple uniform quantization of axis-aligned deltas between consecutive waypoints of agent trajectories.

3. MotionLM

We aim to model a distribution over multi-agent interactions in a general manner that can be applied to distinct downstream tasks, including marginal, joint, and conditional forecasting. This requires an expressive generative

framework capable of capturing the substantial multimodality in driving scenarios. In addition, we take consideration here to preserve temporal dependencies; i.e., inference in our model follows a directed acyclic graph with the parents of every node residing earlier in time and children residing later (Section 3.3, Fig. 4). This enables conditional forecasts that more closely resemble causal interventions [34] by eliminating certain spurious correlations that can otherwise result from disobeying temporal causality². We observe that joint models that do not preserve temporal dependencies may have a limited ability to predict realistic agent reactions – a key use in planning (Section 4.6). To this end, we leverage an autoregressive factorization of our future decoder, where agents’ motion tokens are conditionally dependent on all previously sampled tokens and trajectories are rolled out sequentially (Fig. 2).

Let S represent the input data for a given scenario. This may include context such as roadgraph elements, traffic light states, as well as features describing road agents (e.g., vehicles, cyclists, and pedestrians) and their recent histories, all provided at the current timestep $t = 0$. Our task is to generate predictions for joint agent states $Y_t \doteq \{y_t^1, y_t^2, \dots, y_t^N\}$ for N agents of interest at future timesteps $t = 1, \dots, T$. Rather than complete states, these future state targets are typically two-dimensional waypoints (i.e., (x, y) coordinates), with T waypoints forming the full ground truth trajectory for an individual agent.

²We make no claims that our model is capable of directly modeling causal relationships (due to the theoretical limits of purely observational data and unobserved confounders). Here, we solely take care to avoid breaking temporal causality.

3.1. Joint probabilistic rollouts

In our modeling framework, we sample a predicted action for each target agent at each future timestep. These actions are formulated as discrete motion tokens from a finite vocabulary, as described later in Section 3.2.2. Let a_t^n represent the target action (derived from the ground truth waypoints) for the n th agent at time t , with $A_t \doteq \{a_t^1, a_t^2, \dots, a_t^N\}$ representing the set of target actions for all agents at time t .

Factorization. We factorize the distribution over joint future action sequences as a product of conditionals:

$$p_\theta(A_1, A_2, \dots, A_T | S) = \prod_{t=1}^T p_\theta(A_t | A_{<t}, S), \quad (1)$$

$$p_\theta(A_t | A_{<t}, S) = \prod_{n=1}^N p_\theta(a_t^n | A_{<t}, S). \quad (2)$$

Similar to [36, 43], Eq. (2) represents the fact that we treat agent actions as conditionally independent at time t , given the previous actions and scene context. This aligns empirically with real-world driving over short time intervals; e.g., non-impaired human drivers generally require at least 500 ms to release the accelerator in response to a vehicle braking ahead ([13]). In our experiments, we find 2 Hz reactions to be sufficient to surpass state-of-the-art joint prediction models.

We note that our model’s factorization is entirely latent variable free; multimodal predictions stem purely from categorical token sampling at each rollout timestep.

Training objective. MotionLM is formulated as a generative model trained to match the joint distribution of observed agent behavior. Specifically, we follow a maximum likelihood objective over multi-agent action sequences:

$$\arg \max_{\theta} \prod_{t=1}^T p_\theta(A_t | A_{<t}, S) \quad (3)$$

Similar to the typical training setup of modern language models, we utilize “teacher-forcing” where previous ground truth (not predicted) tokens are provided at each timestep, which tends to increase stability and avoids sampling during training. We note that this applies to all target agents; in training, each target agent is exposed to ground truth action sequence prefixes for all target agents prior to the current timestep. This naturally allows for temporal parallelization when using modern attention-based architectures such as transformers [48].

Our model is subject to the same theoretical limitations as general imitation learning frameworks (e.g., compounding error [38] and self-delusions due to unobserved confounders [33]). However, we find that, in practice, these do not prevent strong performance on forecasting tasks.

3.2. Model implementation

Our model consists of two main networks, an encoder which processes initial scene elements followed by a trajectory decoder which performs both cross-attention to the scene encodings and self-attention along agent motion tokens, following a transformer architecture [48].

3.2.1 Scene encoder

The scene encoder (Fig. 2, left) is tasked with processing information from several input modalities, including the road-graph, traffic light states, and history of surrounding agents’ trajectories. Here, we follow the design of the early fusion network proposed by [31] as the scene encoding backbone of our model. Early fusion is particularly chosen because of its flexibility to process all modalities together with minimal inductive bias.

The features above are extracted with respect to each modeled agent’s frame of reference. Input tensors are then fed to a stack of self-attention layers that exchange information across all past timesteps and agents. In the first layer, latent queries cross-attend to the original inputs in order to reduce the set of vectors being processed to a manageable number, similar to [23, 19]. For additional details, see [31].

3.2.2 Joint trajectory decoder

Our trajectory decoder (Fig. 2, center) is tasked with generating sequences of motion tokens for multiple agents.

Discrete motion tokens. We elect to transform trajectories comprised of continuous waypoints into sequences of discrete tokens. This enables treating sampling purely as a classification task at each timestep, implemented via a standard softmax layer. Discretizing continuous targets in this manner has proven effective in other inherently continuous domains, e.g., in audio generation [46] and mesh generation [30]. We suspect that discrete motion tokens also naturally hide some precision from the model, possibly mitigating compounding error effects that could arise from imperfect continuous value prediction. Likewise, we did not find it necessary to manually add any noise to the ground truth teacher-forced trajectories (e.g., as is done in [26]).

Quantization. To extract target discrete tokens, we begin by normalizing each agent’s ground truth trajectory with respect to the position and heading of the agent at



Figure 3. Displayed are the top two predicted joint rollout modes for three WOMD scenes. Color gradients indicate time progression from $t = 0s$ to $t = 8s$, with the greatest probability joint mode transitioning from green to blue and the secondary joint mode transitioning from orange to purple. Three types of interactions are observed: an agent in the adjacent lane yields to the lane-changing agent according to the timing of the lane change (left), a pedestrian walks behind the passing vehicle according to the progress of the vehicle (center), the turning vehicle either yields to the crossing cyclist (most probable mode) or turns before the cyclist approaches (secondary mode) (right).

time $t = 0$ of the scenario. We then parameterize a uniformly quantized $(\Delta x, \Delta y)$ vocabulary according to a total number of per-coordinate bins as well as maximum and minimum delta values. A continuous, single-coordinate delta action can then be mapped to a corresponding index $\in [0, \text{num_bins} - 1]$, resulting in two indices for a complete $(\Delta x, \Delta y)$ action per step. In order to extract actions that accurately reconstruct an entire trajectory, we employ a greedy search, sequentially selecting the quantized actions that reconstruct the next waypoint coordinates with minimum error.

We wrap the delta actions with a “Verlet” step where a zero action indicates that the same delta index should be used as the previous step (as [36] does for continuous states). As agent velocities tend to change smoothly between consecutive timesteps, this helps reduce the total vocabulary size, simplifying the dynamics of training. Finally, to maintain only T sequential predictions, we collapse the per-coordinate actions to a single integer indexing into their Cartesian product. In practice, for the models presented here, we use 13 tokens per coordinate with $13^2 = 169$ total discrete tokens available in the vocabulary (see Appendix A for further details).

We compute a learned value embedding and two learned positional embeddings (representing the timestep and agent identity) for each discrete motion token, which are combined via an element-wise sum prior to being input to the transformer decoder.

Flattened agent-time self-attention. We elect to include a single self-attention mechanism in the decoder that operates along flattened sequences of all modeled agents’ motion tokens over time. So, given a target sequence of length

T for each of N agents, we perform self-attention over NT elements. While this does mean that these self-attended sequences grow linearly in the number of jointly modeled agents, we note that the absolute sequence length here is still quite small (length 32 for the WOMD interactive split – 8 sec. prediction at 2 Hz for 2 agents). Separate passes of factorized agent and time attention are also possible [32], but we use a single pass here for simplicity.

Ego agent reference frames. To facilitate cross-attention to the agent-centric feature encodings (Section 3.2.1), we represent the flattened token sequence once for each modeled agent. Each modeled agent is treated as the “ego” agent once, and cross-attention is performed on that agent’s scene features. Collapsing the ego agents into the batch dimension allows parallelization during training and inference.

3.3. Enforcing temporal causality

Our autoregressive factorization naturally respects temporal dependencies during joint rollouts; motion token sampling for any particular agent is affected only by past tokens (from any agent) and unaffected by future ones. When training, we require a mask to ensure that the self-attention operation only updates representations at each step accordingly. As shown in Fig. 8 (appendix), this attention mask exhibits a blocked, staircase pattern, exposing all agents to each other’s histories only up to the preceding step.

Temporally causal conditioning. As described earlier, a particular benefit of this factorization is the ability to query for temporally causal conditional rollouts (Fig. 4). In this setting, we fix a query agent to take some sequence of actions and only roll out the other agents.

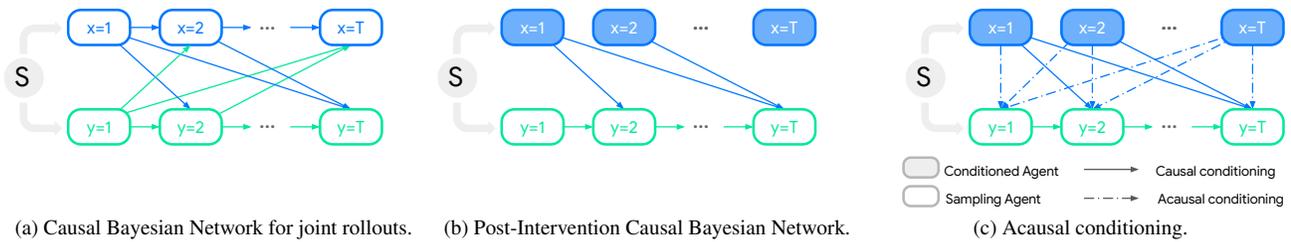


Figure 4. A Causal Bayesian network representation for joint rollouts (left), post-intervention Causal Bayesian network (center), and acausal conditioning (right). Solid lines indicate temporally causal dependencies while dashed lines indicate acausal information flow. Models without temporal dependency constraints will support acausal conditioning but not temporally causal conditioning, which can be problematic when attempting to predict agent reactions.

We may view this as an approximation of computing causal interventions [34] in the absence of confounders; interventions cannot be learned purely through observational data in general (due to the possible presence of unobserved confounders), but our model’s factorization at least eliminates certain spurious correlations arising from breaking temporal causality.

In Fig. 4 (a), we show an example of a Causal Bayesian network governing joint rollouts. Applying an intervention to nodes $x = 1, \dots, T$, by deleting their incoming edges, results in a post-intervention Bayesian network depicted in Fig. 4 (b), which obeys temporal causality. On the other hand, acausal conditioning (Fig. 4 (c)) results in non-causal information flow, where node $x = i$ affects our belief about node $y = j$ for $i \geq j$.

3.4. Rollout aggregation

Joint motion forecasting benchmark tasks like WOMD [14] require a compact representation of the joint future distribution in the form of a small number of joint “modes”. Each mode is assigned a probability and might correspond to a specific homotopic outcome (e.g., pass/yield) or more subtle differences in speed/geometry. Here, we aggregate rollouts to achieve two primary goals: 1) uncover the underlying modes of the distribution and 2) estimate the probability of each mode. Specifically, we follow the non-maximum suppression (NMS) aggregation scheme described in [47], but extend it to the joint setting by ensuring that all agent predictions reside within a given distance threshold to the corresponding cluster. In addition, we leverage model ensembling to account for epistemic uncertainty and further improve the quality of the predictions, combining rollouts from independently trained replicas prior to the aggregation step.

4. Experiments

We evaluate MotionLM on marginal and joint motion forecasting benchmarks, examine its conditional predic-

tions and conduct ablations of our modeling choices.

4.1. Datasets

Waymo Open Motion Dataset (WOMD). WOMD [14] is a collection of 103k 20-second scenarios collected from real-world driving in urban and suburban environments. Segments are divided into 1.1M examples consisting of 9-second windows of interest, where the first second serves as input context and the remaining 8-seconds are the prediction target. Map features such as traffic signal states and lane features are provided along with agent states such as position, velocity, acceleration, and bounding boxes.

Marginal and interactive prediction challenges. For the marginal prediction challenge, six trajectories must be output by the model for each target agent, along with likelihoods of each mode. For the interactive challenge, two interacting agents are labeled in each test example. In this case, the model must output six weighted *joint* trajectories.

4.2. Metrics

The primary evaluation metrics for the marginal and interactive prediction challenges are soft mAP and mAP, respectively, with miss rate as the secondary metric. Distance metrics minADE and minFDE provide additional signal on prediction quality. For the interactive prediction challenge, these metrics refer to scene-level joint calculations. We also use a custom *prediction overlap* metric (similar to [29]) to assess scene-level consistency for joint models. See Appendix C for details on these metrics.

4.3. Model configuration

We experiment with up to 8 model replicas and 512 rollouts per replica, assessing performance at various configurations. For complete action space and model hyperparameter details, see Appendices A and B.

Model	minADE (↓)	minFDE (↓)	Miss Rate (↓)	Soft mAP (↑)
HDGT [20]	0.7676	1.1077	0.1325	0.3709
MPA [22]	0.5913	1.2507	0.1603	0.3930
MTR [40]	0.6050	1.2207	0.1351	0.4216
Wayformer factorized [31]	0.5447	1.1255	0.1229	0.4260
Wayformer multi-axis [31]	0.5454	1.1280	0.1228	0.4335
MTR-A [40]	0.5640	1.1344	0.1160	0.4594
MotionLM (Ours)	0.5509	1.1199	0.1058	0.4507

Table 1. Marginal prediction performance on WOMD test set. We display metrics averaged over time steps (3, 5, and 8 seconds) and agent types (vehicles, pedestrians, and cyclists). Greyed columns indicate the official ranking metrics for the marginal prediction challenge.

Model	minADE (↓)	minFDE (↓)	Miss Rate (↓)	mAP (↑)
SceneTransformer (J) [32]	0.9774	2.1892	0.4942	0.1192
M2I [42]	1.3506	2.8325	0.5538	0.1239
DenseTNT [28]	1.1417	2.4904	0.5350	0.1647
MTR [40]	0.9181	2.0633	0.4411	0.2037
JFP [29]	0.8817	1.9905	0.4233	0.2050
MotionLM (Ours)	0.8911	2.0067	0.4115	0.2178

Table 2. Joint prediction performance on WOMD interactive test set. We display *scene-level* joint metrics averaged over time steps (3, 5, and 8 seconds) and agent types (vehicles, pedestrians, and cyclists). Greyed columns indicate the official ranking metrics for the challenge.

4.4. Quantitative results

Marginal motion prediction. As shown in Table 1, our model is competitive with the state-of-the-art on WOMD marginal motion prediction (independent per agent). For the main ranking metric of soft mAP, our model ranks second, less than 2% behind the score achieved by MTRA [40]. In addition, our model attains a substantially improved miss rate over all prior works, with a relative 9% reduction compared to the previous state-of-the-art. The autoregressive rollouts are able to adequately capture the diversity of multimodal future behavior without reliance on trajectory anchors [8] or static intention points [40].

Interactive motion prediction. Our model achieves state-of-the-art results for the interactive prediction challenge on WOMD, attaining a 6% relative improvement in mAP and 3% relative improvement in miss rate (the two official ranking metrics) over the previous top scoring entry, JFP [29] (see Table 2). In contrast to JFP, our approach does not score pairs of previously constructed marginal trajectories. but generates joint rollouts directly. Fig. 3 displays example interactions predicted by our model.

Table 3 displays prediction overlap rates for various models on the WOMD interactive test and validation sets (see metric details in Appendix C.2). We obtain test set predictions from the authors of [14, 32, 29]. MotionLM obtains the lowest prediction overlap rate, an indication of scene-consistent predictions. In addition, on the validation set we evaluate two versions of our model: marginal and joint. The marginal version does not perform attention

	Model	Prediction Overlap (↓)
Test	LSTM Baseline [14]	0.07462
	Scene Transformer [32]	0.04336
	JFP [29]	0.02671
	MotionLM (joint)	0.02607
Val	MotionLM (marginal)	0.0404
	MotionLM (joint)	0.0292

Table 3. Prediction overlap rates. Displayed is the custom prediction overlap metric for various model configurations on the WOMD interactive test and validation sets.

across the modeled agents during both training and inference rollouts, while the joint version performs 2 Hz interactive attention. We see that the marginal version obtains a relative 38% higher overlap rate than the joint version. The interactive attention in the joint model allows the agents to more appropriately react to one another.

4.5. Ablation studies

Interactive attention frequency. To assess the importance of inter-agent reactivity during the joint rollouts, we vary the frequency of the interactive attention operation while keeping other architecture details constant. For our leaderboard results, we utilize the greatest frequency studied here, 2 Hz. At the low end of the spectrum, 0.125 Hz corresponds to the agents only observing each other’s initial states, and then proceeding with the entire 8-second rollout without communicating again (i.e., marginal rollouts).

Performance metrics generally improve as agents are



Figure 5. Visualization of the top joint rollout mode at the two extremes of the interactive attention frequencies studied here. With no interactive attention (left), the two modeled agents only attend to each other *once* at the beginning of the 8-second rollout and never again, in contrast to 16 total times for 2 Hz attention (right). The independent rollouts resulting from zero interactive attention can result in scene-inconsistent overlap; e.g., a turning vehicle fails to accommodate a crossing pedestrian (top left) or yield appropriately to a crossing vehicle (bottom left).

permitted to interact more frequently (Fig. 6 top, Table 5 in appendix). Greater interactive attention frequencies not only lead to more accurate joint predictions, but also reduce implausible overlaps (i.e., collisions) between different agents’ predictions. Fig. 5 displays examples where the marginal predictions lead to implausible overlap between agents while the joint predictions lead to appropriately separated trajectories. See supplementary for animated visualizations.

Number of rollouts. Our rollout aggregation requires that we generate a sufficient number of samples from the model in order to faithfully represent the multimodal future distribution. For this ablation, we vary the number of rollouts generated, but always cluster down to $k = 6$ modes for evaluation. In general, we see performance metrics improve as additional rollouts are utilized (Fig. 6, bottom and Table 6 in appendix). For our final leaderboard results, we use 512 rollouts per replica, although 32 rollouts is sufficient to surpass the previous top entry on joint mAP.

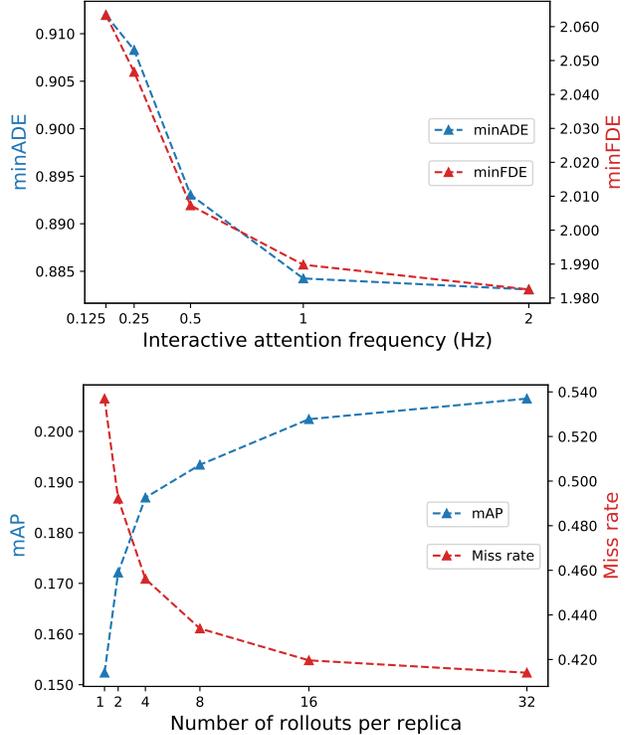


Figure 6. Joint prediction performance across varying interactive attention frequencies (top) and numbers of rollouts per replica (bottom) on the WOMD interactive validation set. Vertical axes display joint (scene-level) metrics for an 8-replica ensemble. See Tables 5 and 6 in the appendix for full parameter ranges and metrics.

4.6. Conditional rollouts

As described in Section 3.3, our model naturally supports “temporally causal” conditioning, similarly to previous autoregressive efforts such as [36, 43]. In this setting, we fix one query agent to follow a specified trajectory and stochastically roll out the target agent. However, we can also modify the model to leak the query agent’s full trajectory, acausally exposing its future to the target agent during conditioning. This resembles the approach to conditional prediction in, e.g., [44], where this acausal conditional distribution is modeled directly, or [29], where this distribution is accessed via inference in an energy-based model.

Here, we assess predictions from our model across three settings: marginal, temporally causal conditional, and acausal conditional (Fig. 4). Quantitatively, we observe that both types of conditioning lead to more accurate predictions for the target agent (Table 4, Fig. 7). Additionally, we see that acausal conditioning leads to greater improvement than temporally causal conditioning relative to marginal predictions across all metrics, e.g., 8.2% increase in soft mAP for acausal vs. 3.7% increase for temporally causal.

Prediction setting	minADE (\downarrow)	minFDE (\downarrow)	Miss Rate (\downarrow)	Soft mAP (\uparrow)
Marginal	0.6069	1.2236	0.1406	0.3951
Temporally causal conditional	0.5997	1.2034	0.1377	0.4096
Acausal conditional	0.5899	1.1804	0.1338	0.4274

Table 4. Conditional prediction performance. Displayed are marginal (single-agent) metrics across three prediction settings for our model on the WOMD interactive validation set: marginal, temporally causal conditional, and acausal conditional.

Intuitively, the greater improvement for acausal conditioning makes sense as it exposes more information to the model. However, the better quantitative scores are largely due to predictions that would be deemed nonsensical if interpreted as predicted *reactions* to the query agent.

This can be illustrated in examples where one agent is following another, where typically the lead agent’s behavior is causing the trailing agent’s reaction, and not vice versa, but this directionality would not be captured with acausal conditioning. This temporally causal modeling is especially important when utilizing the conditional predictions to evaluate safety for an autonomous vehicle’s proposed plans. In a scenario where an autonomous vehicle (AV) is stopped behind another agent, planning to move forward into the other agent’s current position could be viewed as a safe maneuver with acausal conditioning, as the other agent also moving forward is correlated with (but not caused by) the AV proceeding. However, it is typically the lead agent moving forward that causes the trailing AV to proceed, and the AV moving forward on its own would simply rear-end the lead agent.

In the supplementary, we compare examples of predictions in various scenarios for the causal and acausal conditioning schemes. Models that ignore temporal dependencies during conditioning (e.g., [44, 29]) may succumb to the same incorrect reasoning that the acausal version of our model does.

5. Conclusion and future work

In this work, we introduced a method for interactive motion forecasting leveraging multi-agent rollouts over discrete motion tokens, capturing the joint distribution over multimodal futures. The proposed model establishes new state-of-the-art performance on the WOMD interactive prediction challenge.

Avenues for future work include leveraging the trained model in model-based planning frameworks, allowing a search tree to be formed over the multi-agent action rollouts, or learning amortized value functions from large datasets of scene rollouts. In addition, we plan to explore distillation strategies from large autoregressive teachers, enabling faster student models to be deployed in latency-critical settings.



Figure 7. Visualization of the most likely predicted future for the pedestrian in the marginal setting (left) and temporally causal conditional setting (right). When considering the pedestrian independently, the model assigns greatest probability to a trajectory which crosses the road. When conditioned on the the vehicle’s ground truth turn (magenta), the pedestrian is instead predicted to yield.

Acknowledgements. We would like to thank David Weiss, Paul Covington, Ashish Venugopal, Piotr Fidkowski, and Minfa Wang for discussions on autoregressive behavior predictors; Cole Gulino and Brandyn White for advising on interactive modeling; Cheol Park, Wenjie Luo, and Scott Ettinger for assistance with evaluation; Drago Anguelov, Kyracos Shiarlis, and anonymous reviewers for helpful feedback on the manuscript.

References

- [1] Elmira Amirloo, Amir Rasouli, Peter Lakner, Mohsen Rohani, and Jun Luo. LatentFormer: Multi-agent transformer-based interaction modeling and trajectory prediction, 2022.
- [2] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. AudioLM: a language modeling approach to audio generation, 2023.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [4] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *CVPR*, 2020.
- [5] Sergio Casas, Wenjie Luo, and Raquel Urtasun. Intentnet: Learning to predict intention from raw sensor data. In *Conference on Robot Learning*, pages 947–956, 2018.
- [6] Sergio Casas, Cole Gulino, Renjie Liao, and Raquel Urtasun. Spagann: Spatially-aware graph neural networks for relational behavior forecasting from sensor data. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9491–9497. IEEE, 2020.
- [7] Sergio Casas, Cole Gulino, Simon Suo, Katie Luo, Renjie Liao, and Raquel Urtasun. Implicit latent variable model for scene-consistent motion forecasting. In *European Conference on Computer Vision*, pages 624–641. Springer, 2020.
- [8] Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In *Conference on Robot Learning*, pages 86–99. PMLR, 2020.
- [9] Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, et al. Argoverse: 3d tracking and forecasting with rich maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8748–8757, 2019.
- [10] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, et al. PaLM: Scaling language modeling with pathways, 2022.
- [11] Alexander Cui, Sergio Casas, Abbas Sadat, Renjie Liao, and Raquel Urtasun. Lookout: Diverse multi-future prediction and planning for self-driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16107–16116, 2021.
- [12] Henggang Cui, Vladan Radosavljevic, Fang-Chieh Chou, Tsung-Han Lin, Thi Nguyen, Tzu-Kuo Huang, Jeff Schneider, and Nemanja Djuric. Multimodal trajectory predictions for autonomous driving using deep convolutional networks. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 2090–2096. IEEE, 2019.
- [13] Johan Engström, Shu-Yuan Liu, Azadeh Dinparastdjadid, and Camelia Simoiu. Modeling road user response timing in naturalistic settings: a surprise-based framework, 2022. URL <https://arxiv.org/abs/2208.08651>.
- [14] Scott Ettinger, Shuyang Cheng, Benjamin Caine, Chenxi Liu, Hang Zhao, Sabeek Pradhan, Yuning Chai, Benjamin Sapp, Charles R. Qi, Yin Zhou, Zoey Yang, Aurelien Chouard, Pei Sun, Jiquan Ngiam, Vijay Vasudevan, Alexander McCauley, Jonathon Shlens, and Dragomir Anguelov. Large scale interactive motion forecasting for autonomous driving : The waymo open motion dataset. *CoRR*, abs/2104.10133, 2021.
- [15] Roger Girgis, Florian Golemo, Felipe Codevilla, Martin Weiss, Jim Aldon D’Souza, Samira Ebrahimi Kahou, Felix Heide, and Christopher Pal. Latent variable sequential set transformers for joint multi-agent motion prediction. In *International Conference on Learning Representations*, 2021.
- [16] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020.
- [17] Joey Hong, Benjamin Sapp, and James Philbin. Rules of the road: Predicting driving behavior with a convolutional model of semantic interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8454–8462, 2019.
- [18] Maximilian Igl, Daewoo Kim, Alex Kuefler, Paul Mougin, Punit Shah, Kyriacos Shiarlis, Dragomir Anguelov, Mark Palatucci, Brandyn White, and Shimon Whiteson. Symphony: Learning realistic and diverse agents for autonomous driving simulation. In *ICRA*, 2022.
- [19] Andrew Jaegle, Felix Gimeno, Andrew Brock, Andrew Zisserman, Oriol Vinyals, and João Carreira. Perceiver: General perception with iterative attention. In *ICML*, 2021.
- [20] Xiaosong Jia, Penghao Wu, Li Chen, Hongyang Li, Yu Liu, and Junchi Yan. Hdgt: Heterogeneous driving graph transformer for multi-agent trajectory prediction via scene encoding. *CoRL*, 2022.
- [21] Siddhesh Khandelwal, William Qi, Jagjeet Singh, Andrew Hartnett, and Deva Ramanan. What-if motion prediction for autonomous driving. *arXiv preprint arXiv:2008.10587*, 2020.
- [22] Stepan Konev. Mpa: Multipath++ based architecture for motion prediction, 2022. URL <https://arxiv.org/abs/2206.10041>.

- [23] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam R. Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *ICML*, 2019.
- [24] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmohan Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 336–345, 2017.
- [25] Ming Liang, Bin Yang, Rui Hu, Yun Chen, Renjie Liao, Song Feng, and Raquel Urtasun. Learning lane graph representations for motion forecasting. In *European Conference on Computer Vision*, pages 541–556. Springer, 2020.
- [26] Jerry Liu, Wenyuan Zeng, Raquel Urtasun, and Ersin Yumer. Deep structured reactive planning. In *ArXiv*, volume abs/2101.06832, 2021.
- [27] Qiuqing Lu, Weiqiao Han, Jeffrey Ling, Minfa Wang, Haoyu Chen, Balakrishnan Varadarajan, and Paul Covington. Kemp: Keyframe-based hierarchical end-to-end deep model for long-term trajectory prediction. In *2022 International Conference on Robotics and Automation (ICRA)*, 2022.
- [28] Ruikang Luo, Yaofeng Song, Han Zhao, Yicheng Zhang, Yi Zhang, Nanbin Zhao, Liping Huang, and Rong Su. Densent: efficient vehicle type classification neural network using satellite imagery. *ICCV*, 2021.
- [29] Wenjie Luo, Cheol Park, Andre Cornman, Benjamin Sapp, and Dragomir Anguelov. JFP: Structured multi-agent interactive trajectories forecasting for autonomous driving. In *6th Annual Conference on Robot Learning*, 2022.
- [30] Charlie Nash, Yaroslav Ganin, S. M. Ali Eslami, and Peter W. Battaglia. PolyGen: An autoregressive generative model of 3D meshes. *International Conference on Machine Learning*, 2020.
- [31] Nigamaa Nayakanti, Rami Al-Rfou, Aurick Zhou, Kratarth Goel, Khaled S. Refaat, and Benjamin Sapp. Wayformer: Motion forecasting via simple & efficient attention networks. *ArXiv*, abs/2207.05844, 2022.
- [32] Jiquan Ngiam, Vijay Vasudevan, Benjamin Caine, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, David J Weiss, Ben Sapp, Zhifeng Chen, and Jonathon Shlens. Scene transformer: A unified architecture for predicting future trajectories of multiple agents. In *International Conference on Learning Representations*, 2022.
- [33] Pedro A. Ortega, Markus Kunesch, Grégoire Delétang, Tim Genewein, Jordi Grau-Moya, Joel Veness, Jonas Buchli, Jonas Degraeve, Bilal Piot, Julien Pérolat, Tom Everitt, Corentin Tallec, Emilio Parisotto, Tom Erez, Yutian Chen, Scott E. Reed, Marcus Hutter, Nando de Freitas, and Shane Legg. Shaking the foundations: delusions in sequence models for interaction and control. *ArXiv*, abs/2110.10819, 2021.
- [34] Judea Pearl. *Causality*. Cambridge University Press, 2 edition, 2009.
- [35] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [36] Nicholas Rhinehart, Rowan McAllister, Kris Kitani, and Sergey Levine. PRECOG: Prediction conditioned on goals in visual multi-agent settings. In *ICCV*, 2019.
- [37] A Robicquet, A Sadeghian, A Alahi, and S Savarese. Learning social etiquette: Human trajectory prediction in crowded scenes. In *European Conference on Computer Vision (ECCV)*, volume 2, 2020.
- [38] Stephane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*, 2011.
- [39] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Multi-agent generative trajectory forecasting with heterogeneous data for control. In *ECCV*, 2020.
- [40] Shaoshuai Shi, Li Jiang, Dengxin Dai, and Bernt Schiele. Motion transformer with global intention localization and local movement refinement. *Advances in Neural Information Processing Systems*, 2022.
- [41] Haoran Song, Wenchao Ding, Yuxuan Chen, Shaojie Shen, Michael Yu Wang, and Qifeng Chen. PiP: Planning-informed trajectory prediction for autonomous driving. In *ECCV*, 2020.
- [42] Qiao Sun, Xin Huang, Junru Gu, Brian Williams, and Hang Zhao. M2I: From factored marginal trajectory prediction to interactive prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [43] Yichuan Charlie Tang and Ruslan Salakhutdinov. Multiple futures prediction. In *Advances in neural information processing systems*, 2019.
- [44] Ekaterina Tolstaya, Reza Mahjourian, Carlton Downey, Balakrishnan Vadarajan, Benjamin Sapp, and Dragomir Anguelov. Identifying driver interactions via conditional behavior prediction. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- [45] Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International conference on machine learning*, pages 1747–1756. PMLR, 2016.
- [46] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alexander Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. WaveNet: A generative model for raw audio. In *ArXiv*, 2016.

- [47] Balakrishnan Varadarajan, Ahmed Hefny, Avikalp Srivastava, Khaled S. Refaat, Nigamaa Nayakanti, Andre Cornman, Kan Chen, Bertrand Douillard, Chi-Pang Lam, Dragomir Anguelov, and Benjamin Sapp. Multipath++: Efficient information fusion and trajectory aggregation for behavior prediction. In *ICRA*, 2022.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- [49] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*, 2022.
- [50] Ye Yuan, Xinshuo Weng, Yanglan Ou, and Kris M Kitani. Agentformer: Agent-aware transformers for socio-temporal multi-agent forecasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9813–9823, 2021.
- [51] Wei Zhan, Liting Sun, Di Wang, Haojie Shi, Aubrey Clause, Maximilian Naumann, Julius Kummerle, Hendrik Konigshof, Christoph Stiller, Arnaud de La Fortelle, et al. Interaction dataset: An international, adversarial and cooperative motion dataset in interactive driving scenarios with semantic maps. *arXiv preprint arXiv:1910.03088*, 2019.
- [52] Hang Zhao, Jiyang Gao, Tian Lan, Chen Sun, Ben Sapp, Balakrishnan Varadarajan, Yue Shen, Yi Shen, Yuning Chai, Cordelia Schmid, et al. Tnt: Target-driven trajectory prediction. In *Conference on Robot Learning*, pages 895–904. PMLR, 2021.

A. Motion token vocabulary

Delta action space. The models presented in this paper use the following parameters for the discretized delta action space:

- Step frequency: 2 Hz
- Delta interval (per step): [-18.0 m, 18.0 m]
- Number of bins: 128

At 2 Hz prediction, a maximum delta magnitude of 18 m covers axis-aligned speeds up to 36 m/s (~ 80 mph), $> 99\%$ of the WOMD dataset.

Verlet-wrapped action space. Once the above delta action space has the Verlet wrapper applied, we only require 13 bins for each coordinate. This results in a total of $13^2 = 169$ total discrete motion tokens that the model can select from the Cartesian product comprising the final vocabulary.

Sequence lengths. For 8-second futures, the model outputs 16 motion tokens for each agent (note that WOMD evaluates predictions at 2 Hz). For the two-agent interactive split, our flattened agent-time token sequences (Section 3.2.2) have length $2 \times 16 = 32$.

B. Implementation details

B.1. Scene encoder

We follow the design of the early fusion network proposed by [31] as the scene encoding backbone of our model. The following hyperparameters are used:

- Number of layers: 4
- Hidden size: 256
- Feed-forward network intermediate size: 1024
- Number of attention heads: 4
- Number of latent queries: 92
- Activation: ReLU

B.2. Trajectory decoder

To autoregressively decode motion token sequences, we utilize a causal transformer decoder that takes in the motion tokens as queries, and the scene encodings as context. We use the following model hyperparameters:

- Number of layers: 4
- Hidden size: 256



Figure 8. Masked causal attention between two agents during training. We flatten the agent and time axes, leading to an $NT \times NT$ attention mask. The agents may attend to each other’s previous motion tokens (solid squares) but no future tokens (empty squares).

- Feed-forward network intermediate size: 1024
- Number of attention heads: 4
- Activation: ReLU

B.3. Optimization

We train our model to maximize the likelihood of the ground truth motion token sequences via teacher forcing. We use the following training hyperparameters:

- Number of training steps: 600000
- Batch size: 256
- Learning rate schedule: Linear decay
- Initial learning rate: 0.0006
- Final learning rate: 0.0
- Optimizer: AdamW
- Weight decay: 0.6

B.4. Inference

We found nucleus sampling [16], commonly used with language models, to be helpful for improving sample quality while maintaining diversity. Here we set the top- p parameter to 0.95.

C. Metrics descriptions

C.1. WOMD metrics

All metrics for the two WOMD [14] benchmarks are evaluated at three time steps (3, 5, and 8 seconds) and are averaged over all object types to obtain the final value. For joint metrics, a scene is attributed to an object class (vehicle, pedestrian, or cyclist) according to the least common type of agent that is present in that interaction, with cyclist being

the rarest object class and vehicles being the most common. Up to 6 trajectories are produced by the models for each target agent in each scene, which are then used for metric evaluation.

mAP & Soft mAP mAP measures precision of prediction likelihoods and is calculated by first bucketing ground truth futures of objects into eight discrete classes of intent: straight, straight-left, straight-right, left, right, left u-turn, right u-turn, and stationary.

For marginal predictions, a prediction trajectory is considered a “miss” if it exceeds a lateral or longitudinal error threshold at a specified timestep T . Similarly for joint predictions, a prediction is considered a “miss” if none of the k joint predictions contains trajectories for all predicted objects within a given lateral and longitudinal error threshold, with respect to the ground truth trajectories for each agent. Trajectory predictions classified as a miss are labeled as a false positive. In the event of multiple predictions satisfying the miss criteria, consistent with object detection mAP metrics, only one true positive is allowed for each scene, assigned to the highest confidence prediction. All other predictions for the object are assigned a false positive.

To compute the mAP metric, bucket entries are sorted and a P/R curve is computed for each bucket, averaging precision values over various likelihood thresholds for all intent buckets results in the final mAP value. Soft mAP differs only in the fact that additional matching predictions (other than the most likely match) are ignored instead of being assigned a false positive, and so are not penalized in the metric computation.

Miss rate Using the same definition of a “miss” described above for either marginal or joint predictions, miss rate is a measure of what fraction of scenarios fail to generate *any* predictions within the lateral and longitudinal error thresholds, relative to the ground truth future.

minADE & minFDE minADE measures the Euclidean distance error averaged over all timesteps for the closest prediction, relative to ground truth. In contrast, minFDE considers only the distance error at the final timestep. For joint predictions, minADE and minFDE are calculated as the average value over both agents.

C.2. Prediction overlap

As described in [29], the WOMD [14] overlap metric only considers overlap between predictions and ground truth. Here we use a *prediction overlap* metric to assess scene-level consistency for joint models. Our implementation is similar to [29], except we follow the convention of the WOMD challenge of only requiring models to generate (x, y) waypoints; headings are inferred as in [14]. If

the bounding boxes of two predicted agents collide at any timestep in a scene, that counts as an overlap/collision for that scene. The final prediction overlap rate is calculated as the sum of per-scene overlaps, averaged across the dataset.

D. Additional evaluation

Ablations. Tables 5 and 6 display joint prediction performance across varying interactive attention frequencies and numbers of rollouts, respectively. In addition to the ensemble model performance, single replica performance is evaluated. Standard deviations are computed for each metric over 8 independently trained replicas.

Scaling analysis. Table 7 displays the performance of different model sizes on the WOMD interactive split, all trained with the same optimization hyperparameters. We vary the number of layers, hidden size, and number of attention heads in the encoder and decoder proportionally. Due to external constraints, in this study we only train a single replica for each parameter count. We observe that a model with 27M parameters overfits while 300K underfits. Both the 1M and 9M models perform decently. In this paper, our main results use 9M-parameter replicas.

Latency analysis. Table 8 provides inference latency on the latest generation of GPUs across different numbers of rollouts. These were measured for a single-replica joint model rolling out two agents.

E. Visualizations

In the supplementary zip file, we have included GIF animations of the model’s greatest-probability predictions in various scenes. Each example below displays the associated scene ID, which is also contained in the corresponding GIF filename. We describe the examples here.

E.1. Marginal vs. Joint

- Scene ID: 286a65c777726df3
Marginal: The turning vehicle and crossing cyclist collide.
Joint: The vehicle yields to the cyclist before turning.
- Scene ID: 440bbf422d08f4c0
Marginal: The turning vehicle collides with the crossing vehicle in the middle of the intersection.
Joint: The turning vehicle yields and collision is avoided.
- Scene ID: 38899bce1e306fb1
Marginal: The lane-changing vehicle gets rear-ended by the vehicle in the adjacent lane.
Joint: The adjacent vehicle slows down to allow the lane-changing vehicle to complete the maneuver.

Freq. (Hz)	Ensemble				Single Replica			
	minADE (\downarrow)	minFDE (\downarrow)	MR (\downarrow)	mAP (\uparrow)	minADE (\downarrow)	minFDE (\downarrow)	MR (\downarrow)	mAP (\uparrow)
0.125	0.9120	2.0634	0.4222	0.2007	1.0681 (0.011)	2.4783 (0.025)	0.5112 (0.007)	0.1558 (0.007)
0.25	0.9083	2.0466	0.4241	0.1983	1.0630 (0.009)	2.4510 (0.025)	0.5094 (0.006)	0.1551 (0.006)
0.5	0.8931	2.0073	0.4173	0.2077	1.0512 (0.009)	2.4263 (0.022)	0.5039 (0.006)	0.1588 (0.004)
1	0.8842	1.9898	0.4117	0.2040	1.0419 (0.014)	2.4062 (0.032)	0.5005 (0.008)	0.1639 (0.005)
2	0.8831	1.9825	0.4092	0.2150	1.0345 (0.012)	2.3886 (0.031)	0.4943 (0.006)	0.1687 (0.004)

Table 5. Joint prediction performance across varying interactive attention frequencies on the WOMD interactive validation set. Displayed are *scene-level* joint evaluation metrics. For the single replica metrics, we include the standard deviation (across 8 replicas) in parentheses.

# Rollouts	Ensemble				Single Replica			
	minADE (\downarrow)	minFDE (\downarrow)	MR (\downarrow)	mAP (\uparrow)	minADE (\downarrow)	minFDE (\downarrow)	MR (\downarrow)	mAP (\uparrow)
1	1.0534	2.3526	0.5370	0.1524	1.9827 (0.018)	4.7958 (0.054)	0.8182 (0.003)	0.0578 (0.004)
2	0.9952	2.2172	0.4921	0.1721	1.6142 (0.011)	3.8479 (0.032)	0.7410 (0.003)	0.0827 (0.004)
4	0.9449	2.1100	0.4561	0.1869	1.3655 (0.012)	3.2060 (0.035)	0.6671 (0.003)	0.1083 (0.003)
8	0.9158	2.0495	0.4339	0.1934	1.2039 (0.013)	2.7848 (0.035)	0.5994 (0.004)	0.1324 (0.003)
16	0.9010	2.0163	0.4196	0.2024	1.1254 (0.012)	2.5893 (0.031)	0.5555 (0.005)	0.1457 (0.003)
32	0.8940	2.0041	0.4141	0.2065	1.0837 (0.013)	2.4945 (0.035)	0.5272 (0.005)	0.1538 (0.004)
64	0.8881	1.9888	0.4095	0.2051	1.0585 (0.012)	2.4411 (0.033)	0.5114 (0.005)	0.1585 (0.004)
128	0.8851	1.9893	0.4103	0.2074	1.0456 (0.012)	2.4131 (0.033)	0.5020 (0.006)	0.1625 (0.004)
256	0.8856	1.9893	0.4078	0.2137	1.0385 (0.012)	2.3984 (0.031)	0.4972 (0.007)	0.1663 (0.005)
512	0.8831	1.9825	0.4092	0.2150	1.0345 (0.012)	2.3886 (0.031)	0.4943 (0.006)	0.1687 (0.004)

Table 6. Joint prediction performance across varying numbers of rollouts per replica on the WOMD interactive validation set. Displayed are *scene-level* joint evaluation metrics. For the single replica metrics, we include the standard deviation (across 8 replicas) in parentheses.

Parameter count	Miss Rate (\downarrow)	mAP (\uparrow)
300K	0.6047	0.1054
1M	0.5037	0.1713
9M	0.4972	0.1663
27M	0.6072	0.1376

Table 7. Joint prediction performance across varying model sizes on the WOMD interactive validation set. Displayed are *scene-level* joint mAP and miss rate for 256 rollouts for a single model replica (except for 9M which displays the mean performance of 8 replicas).

- Scene ID: 2ea76e74b5025ec7
Marginal: The cyclist crosses in front of the vehicle leading to a collision.
Joint: The cyclist waits for the vehicle to proceed before turning.
- Scene ID: 55b5fe989aa4644b
Marginal: The cyclist lane changes in front of the adjacent vehicle, leading to collision.
Joint: The cyclist remains in their lane for the duration of the scene, avoiding collision.

Number of rollouts	Latency (ms)
16	19.9 (0.19)
32	27.5 (0.25)
64	43.8 (0.26)
128	75.8 (0.23)
256	137.7 (0.19)

Table 8. Inference latency on current generation of GPUs for different numbers of rollouts of the joint model. We display the mean and standard deviation (in parentheses) of the latency measurements for each setting.

E.2. Marginal vs. Conditional

“Conditional” here refers to temporally causal conditioning as described in the main text.

- Scene ID: 5ebba77f351358e2
Marginal: The pedestrian crosses the street as a vehicle is turning, leading to a collision.
Conditional: When conditioning on the vehicle’s turning trajectory as a query, the pedestrian is instead predicted to remain stationary.
- Scene ID: d557eee96705c822

Marginal: The modeled vehicle collides with the lead vehicle.

Conditional: When conditioning on the lead vehicle's query trajectory, which remains stationary for a bit, the modeled vehicle instead comes to an appropriate stop.

- Scene ID: 9410e72c551f0aec

Marginal: The modeled vehicle takes the turn slowly, unaware of the last turning vehicle's progress.

Conditional: When conditioning on the query vehicle's turn progress, the modeled agent likewise makes more progress.

- Scene ID: c204982298bd1a1

Marginal: The modeled vehicle proceeds slowly, unaware of the merging vehicle's progress.

Conditional: When conditioning on the query vehicle's merge progress, the modeled agent accelerates behind.

E.3. Temporally Causal vs. Acausal Conditioning

- Scene ID: 4f39d4eb35a4c07c

Joint prediction: The two modeled vehicles maintain speed for the duration of the scene.

Conditioning on trailing agent:

- **Temporally causal:** The lead vehicle is indifferent to the query trailing vehicle decelerating to a stop, proceeding along at a constant speed.

- **Acausal:** The lead vehicle is "influenced" by the query vehicle decelerating. It likewise comes to a stop. Intuitively, this is an incorrect direction of influence that the acausal model has learned.

Conditioning on lead agent:

- **Temporally causal:** When conditioning on the query lead vehicle decelerating to a stop, the modeled trailing vehicle is likewise predicted to stop.

- **Acausal:** In this case, the acausal conditional prediction is similar to the temporally causal conditional. The trailing vehicle is predicted to stop behind the query lead vehicle.