

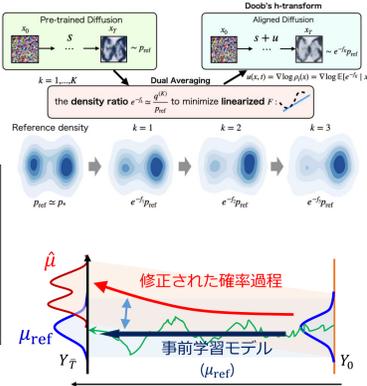
生成モデルの新しいアライメント手法

拡散モデルのアライメント [Kawata, Oko, Nitanda, Suzuki: Direct Distributional Optimization for Provable Alignment of Diffusion Models. ICLR2025]

拡散モデルの生成分布を直接修正することはこれまで困難だった。⇒本研究: 直接アライメントを修正する手法を提案

$$\min_{\mu \in \mathcal{P}} \mathcal{L}(\mu) = F(\mu) + \beta \text{KL}(\mu || \mu_{\text{ref}})$$

例: DPO ベイズフィルタ



Phase 1: 最適化. 双対平均化

For $k = 1, \dots, N-1$:

$$\mu^{(k+1)} = \arg \min_{\mu \in \mathcal{P}} \frac{2}{k(k+1)} \sum_{j=1}^k \left(\mathbb{E}_{\mu} \left[\frac{\delta F(\mu^{(j)})}{\delta \mu} \right] + \beta \text{KL}(\mu || \mu_{\text{ref}}) \right) + \frac{2\beta}{k} \text{KL}(\mu || \mu_{\text{ref}})$$

where $\hat{g}^{(k)} = \sum_{j=1}^k \frac{j}{\beta(k+1)(k+2)/2} \frac{\delta F(\mu^{(j)})}{\delta \mu}$ ← ニューラルネット近似

→ $\frac{d\hat{\mu}}{d\mu_{\text{ref}}} \propto \exp(-\hat{g})$ μ_{ref} と最適解 $\hat{\mu}$ の密度比
ただし $\hat{g} \leftarrow \hat{g}^{(N)}$ のみが見られる。

Phase 2: Sampling from $\hat{\mu} \propto \exp(-\hat{g}) \mu_{\text{ref}}$.

Doob h-Transform (Doob, 1957; Rogers & Williams, 2000)

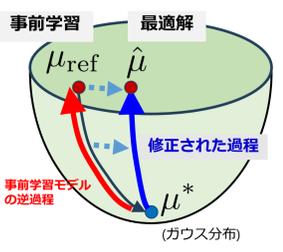
$$d\tilde{Y}_t = (\tilde{Y}_t + 2s(\tilde{Y}_t, \bar{T} - t) + 2\nabla_x \log(\mathbb{E}[\exp(-\hat{g}(Y_T)) | Y_t = x])|_{x=\tilde{Y}_t})dt + \sqrt{2}dW_t$$

修正項

$$\tilde{Y}_T \sim \exp(-\hat{g})\mu_{\text{ref}} = \hat{\mu}$$

元の μ_{ref} の逆過程: $dY_t = (Y_t + 2s(Y_t, \bar{T} - t))dt + \sqrt{2}dW_t$

最適解 $\hat{\mu}$ からのサンプリングが可能に!



LLMのアライメント 人間の嗜好モデルを特別に仮定せずアライメントさせたい。

各プロンプト $c \sim p(c)$ に対し以下のデータを得る:

$y_1 \sim p_w(y|c), y_2 \sim p_l(y|c)$

※ y_1, y_2 は必ずしも対になっていない。

仮定: $p_{\text{ref}} = \theta p_w + (1 - \theta)p_l$

目標: p_w からデータ生成できるようにしたい。
→ 特定の嗜好モデルに依存しない。統計的整合性を保証したい。

提案法「Density ratio matching estimator」:

真の密度比からの距離 KL-正則化

$$\min_{p: \text{LLM}} \{ B_f(r || r^*) + \beta^{-1} \text{KL}(p || p_{\text{ref}}) \}$$

ただし, $r^* = \frac{p_l}{p_w}$ かつ $r = (1 - \theta)^{-1} \left(\frac{p_{\text{ref}}}{p} - \theta \right)$.

($r = r^* \Leftrightarrow p = p_w$)

$B_f(r || r^*) := \int [f(r^*) - f(r) - f'(r)(r^* - r)] dp_w$

: 凸関数 f で決まる Bregman ダイバージェンス

[Higuchi, Suzuki: Direct Density Ratio Optimization: A Statistically Consistent Approach to Aligning Large Language Models. ICML2025]

$$B_f(r || r^*) + \beta^{-1} \text{KL}(p || p_{\text{ref}})$$

$$= \int (-f(r) + f'(r)r) dp_w - \int f'(r) dp_l + \beta^{-1} \int \hat{r}^{-1} \log(\hat{r}) dp_{\text{ref}}$$

→ データから推定可能な形

→ $r = (1 - \theta)^{-1} \left(\frac{p_{\text{ref}}}{p} - \theta \right)$ と $\hat{r} = \frac{p_{\text{ref}}}{p}$ を代入して, p に関して最適化

例 (logistic regression):

$$f(t) = t \log(t) - (1+t) \log(1+t)$$

$$\int \log(1+r) dp_w + \int \log(1+r^{-1}) dp_l + \beta^{-1} \int \hat{r}^{-1} \log(\hat{r}) dp_{\text{ref}}$$

Logistic損失 KL-正則化

統計的整合性の理論保証有り
→ 「相対的密度比」を用いることで安定化可能



線形注意手法の理論と最適状態サイズの決定

• 現状のAttentionは本当にそれで十分か?
• 計算量の問題がある。

Attention	<	代替手法 (線形手法)
• 計算量 $O(L^2)$ (L は入力長)	効率性	• 計算量 $O(L)$ or $O(L \log(L))$
• 基盤モデルのスタンダード	能力	• 能力の限界に関する指摘 言語等の離散データが不得意

Q. 効率的な代替法はあるか?
• どのようなタスクなら代替可能?

状態空間モデルの表現力 [Naoki Nishikawa, Taiji Suzuki: State Space Models are Comparable to Transformers in Estimating Functions with Dynamic Smoothness. ICLR2025]

• **主結果 1** 既存結果: Copying タスクで 1層のSSM は Transformer を代替できない [Jelassi et al.: Repeat After Me: Transformers are Better than State Space Models at Copying. 2024]



なぜか? ⇒ 入力依存で重要なトークンを抽出可能する必要はある

本研究: 多層の FNN + SSM で Transformer を代替可能
各トークンの重要度を前の層で計算すれば実は代替できる:

p 1 u 7 v 4 w 7 t 9 u ?
1 3 20 159 2 2 4 3 2 3 24

← 重要度 (入力依存) e.g. 自分 or 1つ前の最後のトークンと同じ

• **主結果 2** 区分的 γ -平滑関数の推定において, SSMとTransformer は同じ推定誤差を達成することを証明

線形注意機構の特徴次元の決定 [Naoki Nishikawa, Rei Higuchi, Taiji Suzuki: Degrees of Freedom for Linear Attention: Distilling Softmax Attention with Optimal Feature Efficiency. NeurIPS2025]

• Attentionのカーネル関数としての定式化と分解:
 $K(x, y) = \exp(x^T y / \sqrt{d}) = \mathbb{E}_{z \sim N(0, I)} [\phi(x; z) \phi(y; z)]$ where $\phi(x; z) = \exp\left(\frac{z^T x}{d^{1/4}} - \frac{\|x\|^2}{2\sqrt{d}}\right)$

• 線形注意機構と特徴写像: $\hat{K}(x, y) = \varphi(x)^T \varphi(y)$ where $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^M$.

定理 (近似誤差の理論評価)

$\lambda > 0$ に対し, N_λ を Attention の統計的自由度とする。すると,

$$M \geq \frac{4}{\lambda} N_\lambda \log\left(\frac{32N_\lambda}{t/2}\right)$$

とすることで, Attentionは線形注意機構で以下の誤差で近似できる:

$$\|K - \hat{K}\|_{L^2(P_X \otimes P_X)}^2 \leq 2t \left(\lambda^2 + \|K\|_{L^2(P_X \otimes P_X)}^2 \right)$$

(ただし, 特徴写像 φ は適切に選ぶとする)

• Models: GPT-2 (1.24 × 10⁸ parameters), Pythia-1B (1.01 × 10⁹ parameters) [Radford et al., 2019] [Biderman et al., 2023]

• Hyperparameters: (C, λ) = (64, 10⁻⁴) for GPT-2, (C, λ) = (128, 10⁻⁸) for Pythia-1B

✓ Cost C is set to be the same as the head size following prior work (e.g. [Chen et al., 2025])

Model	Cost	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
GPT-2	64	130	182	35	44	42	65	92	33	28	34	46	39	-	-	-	-
Pythia-1B	128	277	138	275	132	204	167	115	142	231	64	96	73	40	52	34	9

Large in early and middle layers Small in latter layers
深いレイヤーの方が少ない特徴量でOK (CNN, FNNと同様の現象 [Arora et al., 2018; Ravichandran et al., 2019; Suzuki et al., 2020])

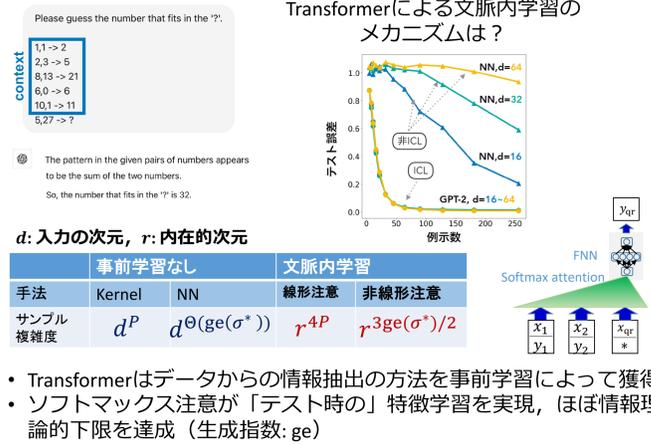
Downstream taskの性能 対抗手法: Performer [Choromanski et al., 2021], DiJiang [Chen et al., 2025]

Strategy	Method	PiQA	logiQA	ARC-E	ARC-C	Winogrande	MMLU	WSC	Average
Original GPT-2	direct	0.5985	0.3103	0.3325	0.3003	0.5122	0.2789	0.6538	0.4266
	DiJiang	0.5065	0.2550	0.2113	0.2244	0.4846	0.2639	0.4615	0.3409
Fix	Performer	0.5468	0.2934	0.3039	0.2747	0.4996	0.2517	0.5962	0.3952
	direct	0.5832	0.3195	0.2921	0.2995	0.5335	0.2552	0.6154	0.4141
	softmax	0.5718	0.2673	0.2479	0.3029	0.5020	0.2634	0.6154	0.3958
DoF	L^2	0.5822	0.3195	0.2483	0.2773	0.5107	0.2520	0.5962	0.3980
	direct	0.5669	0.3011	0.3241	0.3012	0.5280	0.2712	0.6346	0.4182
	softmax	0.5751	0.3026	0.3224	0.2995	0.5328	0.2564	0.6442	0.4190
DoF + Clip	L^2	0.5664	0.3088	0.2736	0.2824	0.4972	0.2608	0.5865	0.3965
	direct	0.5892	0.3164	0.3136	0.2952	0.5075	0.2993	0.6346	0.4223
	softmax	0.5860	0.3026	0.3401	0.2816	0.4996	0.2832	0.6346	0.4182
L^2	0.5822	0.3164	0.2942	0.3063	0.5091	0.2799	0.5673	0.4079	

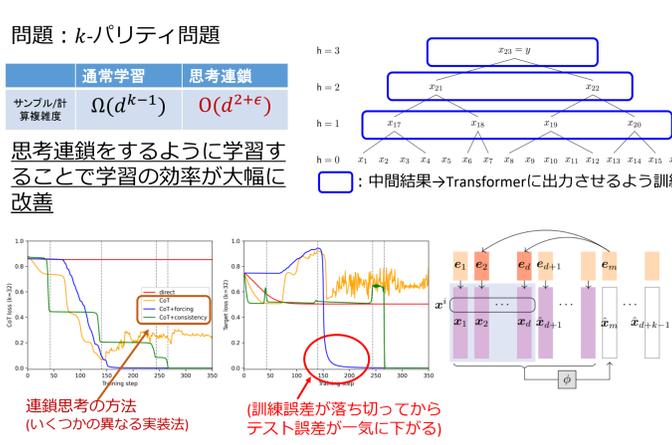
(特徴次元を固定した方法) (特徴次元を層ごとに選んだ方法) (提案法)

テスト時推論の理論

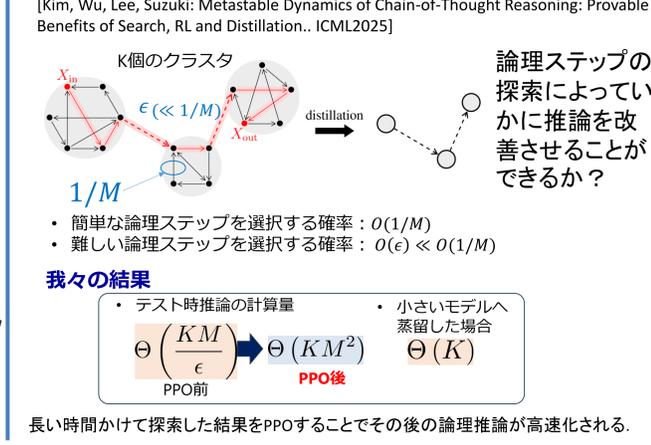
文脈内学習 [Nishikawa, Song, Oko, Wu, Suzuki: Nonlinear transformers can perform inference-time feature learning. ICML2025]



思考連鎖 [Kim, Suzuki: Transformers Provably Solve Parity Efficiently with Chain of Thought. ICLR2025]

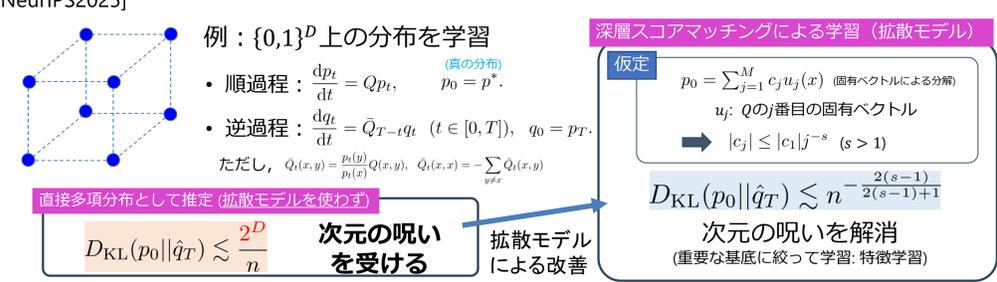


テスト時スケールアップ [Kim, Wu, Lee, Suzuki: Metastable Dynamics of Chain-of-Thought Reasoning: Provable Benefits of Search, RL and Distillation. ICML2025]



離散拡散モデルの理論

[Wakasugi, Suzuki: State Size Independent Statistical Error Bound for Discrete Diffusion Models. NeurIPS2025]



確率測度の非凸目的関数の最適化

[Naoya Yamamoto, Juno Kim, Taiji Suzuki: Hessian-guided Perturbed Wasserstein Gradient Flows for Escaping Saddle Points. NeurIPS2025]

