



Modern Recursive Transformer

Taekhyun Park

Department of Data Science,
Pusan National University, Busan
pthpark1@pusan.ac.kr

Contents

■ Introduction

- Recursive Deep Learning Model
- How to make good Recursive System for Transformer

■ Modern Recursive Model

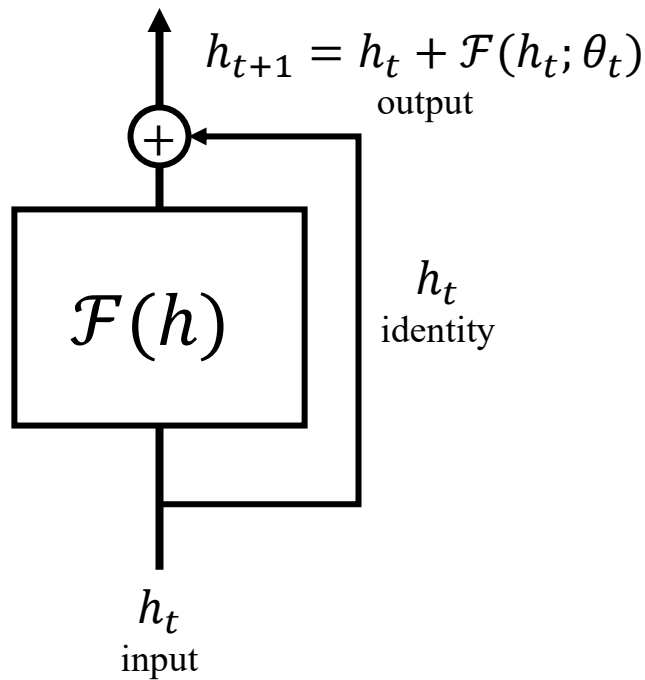
- Autoregressive based Transformer
- Diffusion based Transformer
- Energy based Transformer
- Hierarchical Reasoning Model
- Scaling Latent Reasoning

■ Future Work

- Deep Up Scaling
- Future Work

Recursive Deep Learning Model

- Since the advent of ResNet (He et al., 2015), deep learning models can be interpreted as **continuous-time dynamical systems**. As training aligns hidden representations within a shared latent space with a basis \mathcal{B} , **recursion design** becomes critical to performance and stability (Liu et al., 2026).



Discrete formulation :

$$h_{t+1} = h_t + \mathcal{F}_{\theta_t}(h_t)$$

Continuous formulation :

$$\frac{dh(t)}{dt} = F_{\theta}(h(t), t)$$

How to make good Recursive systems for Transformer

1. Scaling of Transformer

- Human cognition is described as two systems: **System 1**, which enables fast(30~100Hz, fast gamma waves), **intuitive thinking** and **System2**, which supports slow(4~8Hz, slow theta waves), **deliberate reasoning** (Kahneman, 2011).

SYSTEM 1

Intuition & instinct

95%

Unconscious
Fast
Associative
Automatic pilot

SYSTEM 2

Rational thinking

5%

Takes effort
Slow
Logical
Lazy
Indecisive

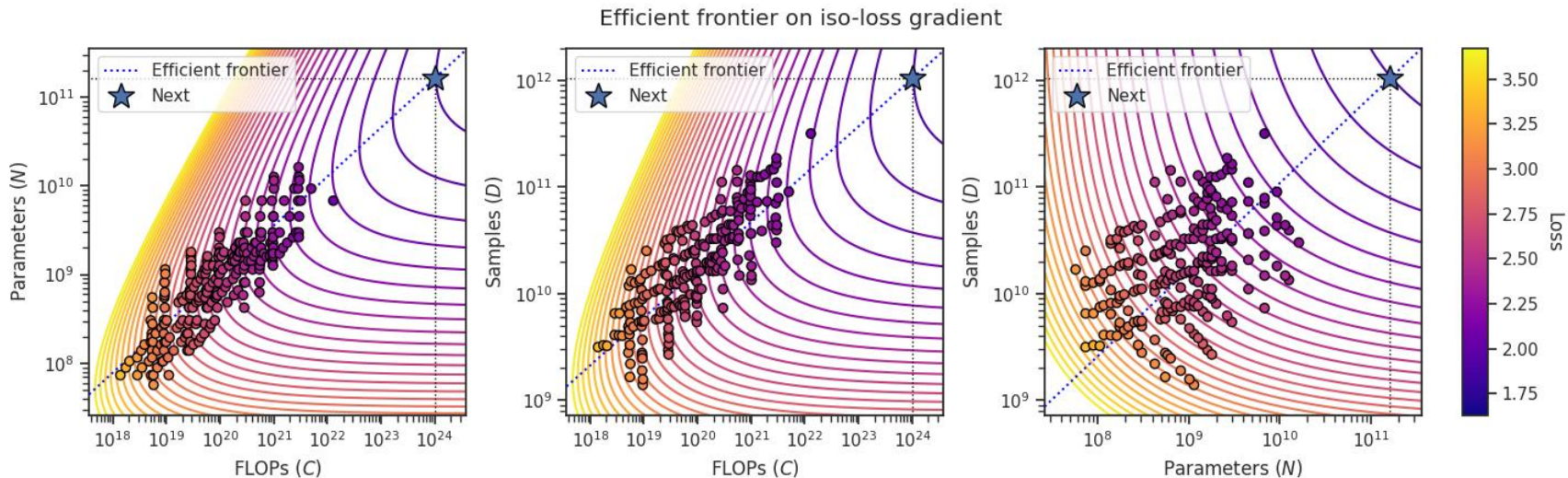


Source: Daniel Kahneman

How to make good Recursive systems for Transformer

1. Scaling of Transformer

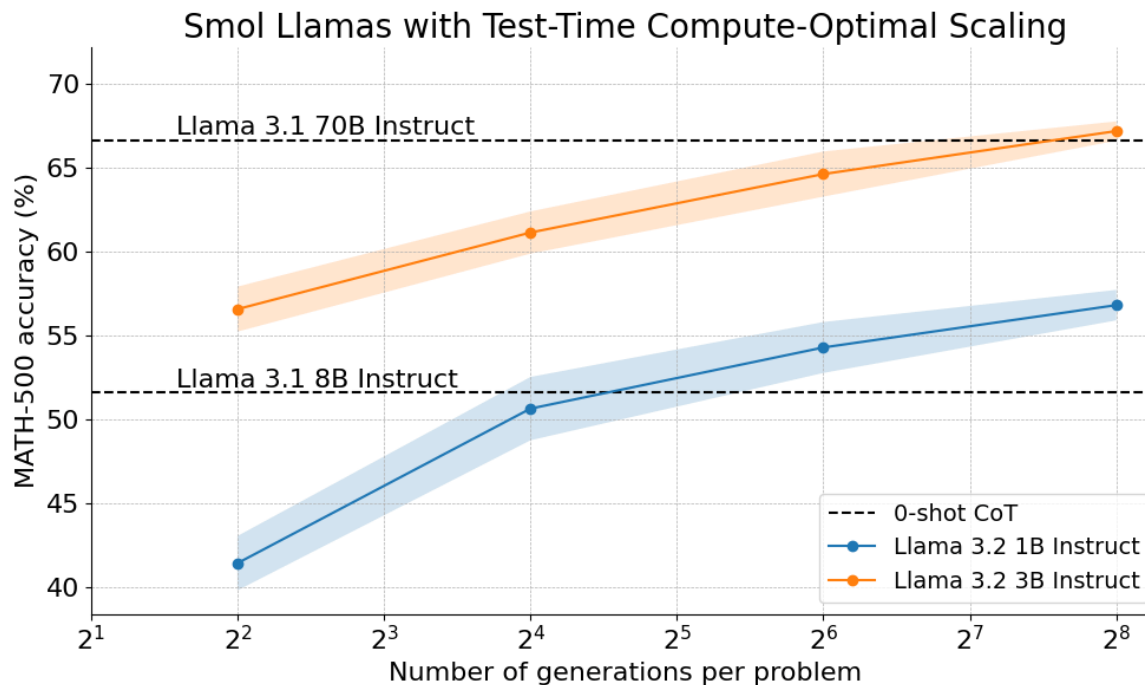
- Inspired by this distinction, Transformer research emphasizes reasoning depth, either by increasing **Test-Time Compute(TTC)**, spending more computation per token (Snell et al., 2025) and large-scale training such as **Chinchilla scaling** (Hoffmann et al., 2022).
- Data scarcity limits further scaling, and recent studies show that **TTC** effectively increase recursive reasoning depth, with explicit mechanisms such as **CoT** (Wei et al., 2022) yielding stronger performance gains beyond a certain scale(Snell et al., 2025).



How to make good Recursive systems for Transformer

1. Scaling of Transformer

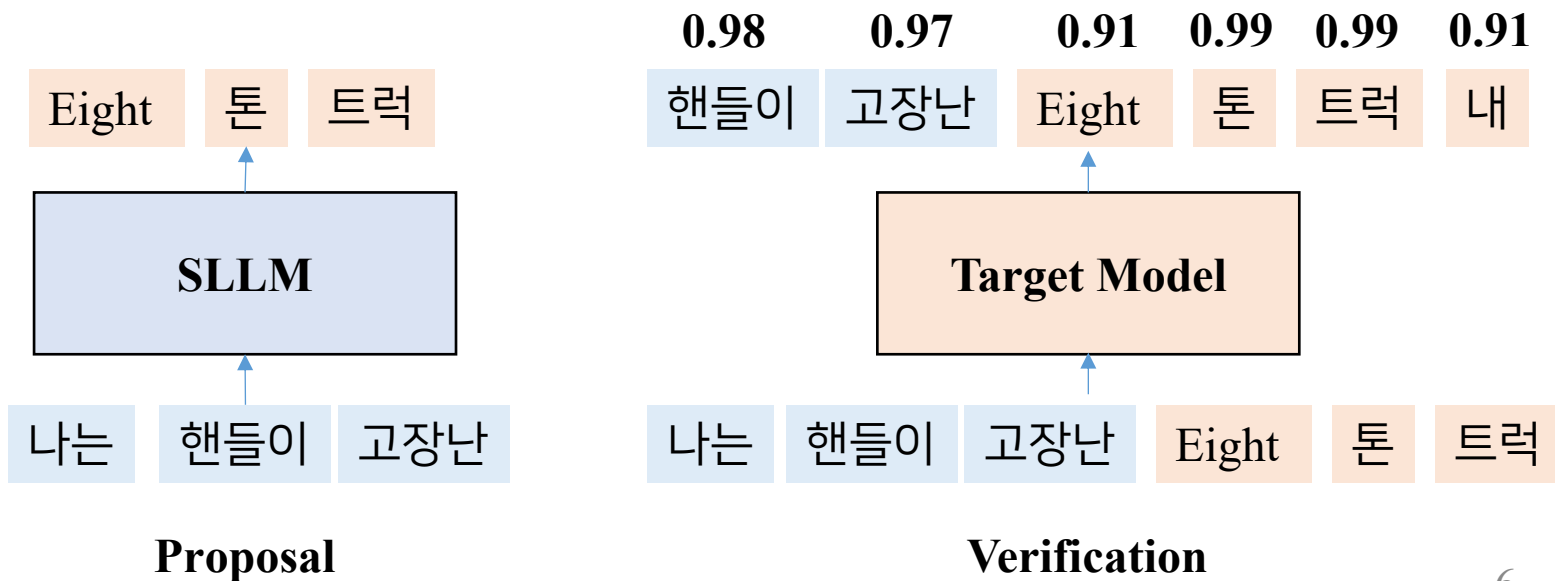
- Inspired by this distinction, Transformer research emphasizes reasoning depth, either by increasing **Test-Time Compute**(TTC), spending more computation per token (Snell et al., 2025) or by large-scale training such as **Chinchilla scaling** (Hoffmann et al., 2022).
- Data scarcity limits further scaling, and recent studies show that **TTC** effectively increase recursive reasoning depth, with explicit mechanisms such as **CoT** (Wei et al., 2022) yielding stronger performance gains beyond a certain scale (Snell et al., 2025).



How to make good Recursive systems for Transformer

2. Scalable Learning

- However, the performance gains achieved through TTC and Chinchilla-style scaling come with clear trade-offs in computational cost and inference latency.
 - **Speculative Decoding**(Leviathan et al., 2023), first generates candidate outputs with a smaller model and then verifies them with a larger model; if verification fails, the large model recomputes the corresponding parts. (**Google Gemini 3.0**)
 - **Routing-based methods** use lightweight classifiers to route easy queries to smaller models and harder queries to more powerful ones. (**OpenAI GPT5**)
 - **Automated approaches** enable models to dynamically scale their own recursion depth during inference, as seen in architectures such as **Hierarchical Reasoning Model (HRM)**(Wang et al., 2025) and **Energy-based Transformer**(Gladstone et al., 2026).

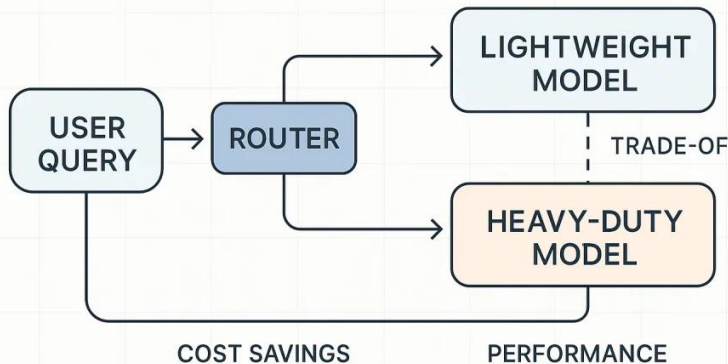


How to make good Recursive systems for Transformer

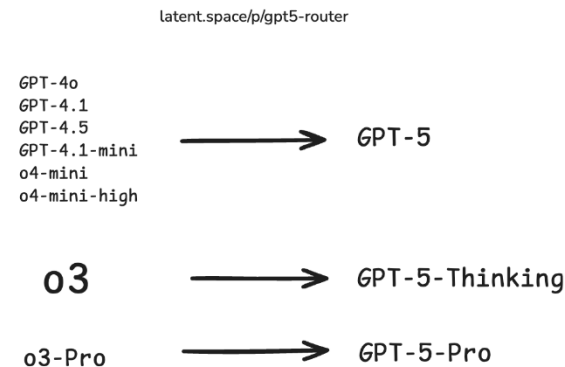
2. Scalable Learning

- However, the performance gains achieved through TTC and Chinchilla-style scaling come with clear trade-offs in computational cost and inference latency.
 - **Speculative Decoding**(Leviathan et al., 2023), first generates candidate outputs with a smaller model and then verifies them with a larger model; if verification fails, the large model recomputes the corresponding parts. (**Google Gemini 3.0**)
 - **Routing-based methods** use lightweight classifiers to route easy queries to smaller models and harder queries to more powerful ones. (**OpenAI GPT5**)
 - **Automated approaches** enable models to dynamically scale their own recursion depth during inference, as seen in architectures such as **Hierarchical Reasoning Model (HRM)**(Wang et al., 2025) and **Energy-based Transformer**(Gladstone et al., 2026).

Inside GPT-5: The Dual-Model System Explained



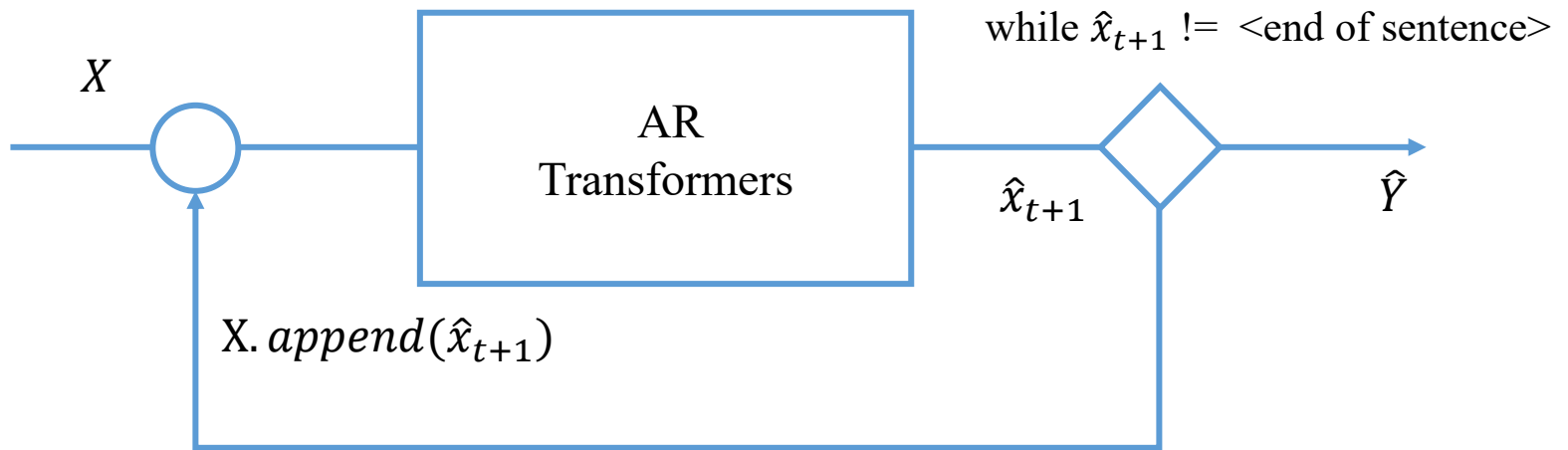
Mapping Deprecated models to GPT-5 equivalents (for ChatGPT)



Autoregressive based Transformer

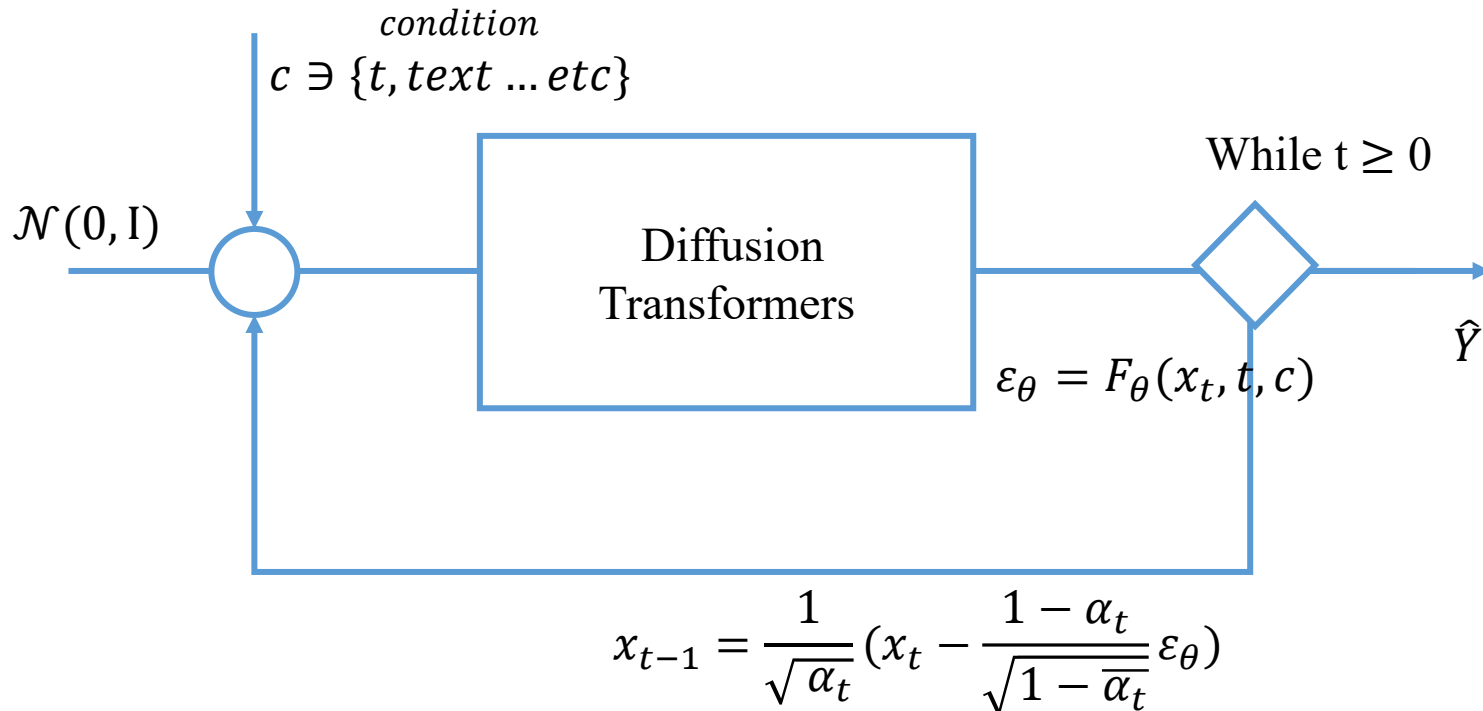
- To expand TTC, GRPO(Shao, Zhihong, et al., 2024) is used to induce CoT reasoning, leading to improved performance.

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]$$
$$\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min \left[\frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,<t})} \hat{A}_{i,t}, \text{clip} \left(\frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{KL} [\pi_{\theta} || \pi_{ref}] \right\}, \quad (3)$$



Diffusion based Transformer

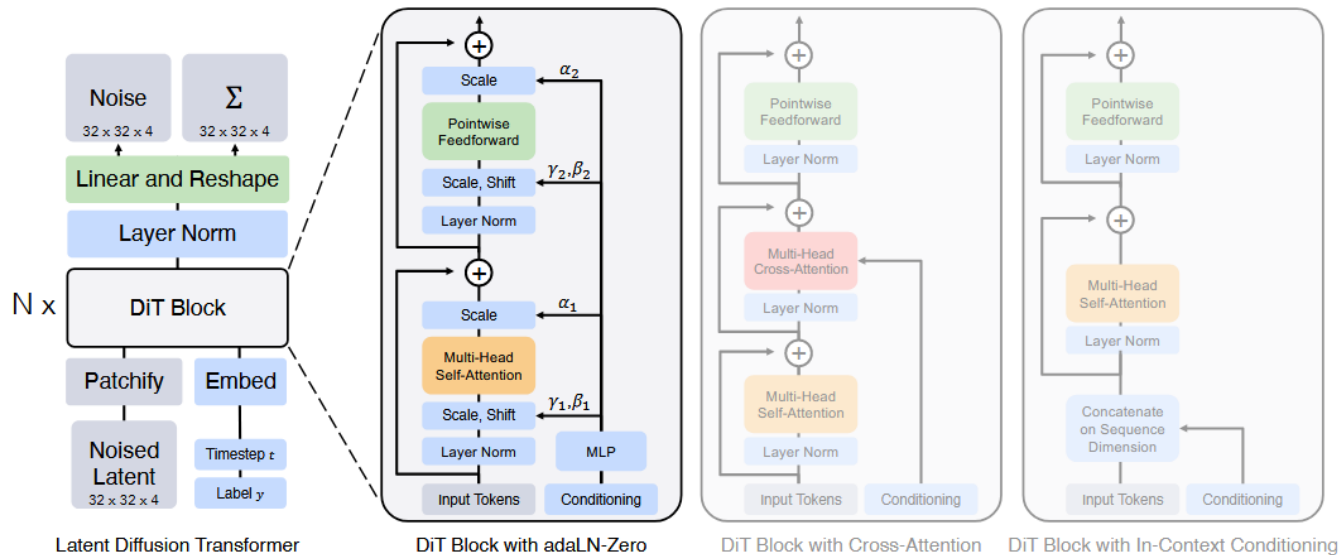
- Diffusion Transformer(DiT)(Peebles et al., 2023, Nie et al, Neurips 2025) defines generation as a diffusion process based on Brownian-motion-driven SDEs , where the model **estimates the assumed noise component ε_θ** and **iteratively removes** it according to a predefined noise schedule at each step to recover clean representations.



Diffusion based Transformer

Scalable Diffusion Models with Transformers (ICCV 2023) (Peebles et al., 2023)

- Forward process: $x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon_t$, where $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$
- Loss (training): $\mathcal{L}(\theta) = -p(x_0|x_1) + \sum_t D_{KL}(q^*(x_{t-1}|x_t, x_0)|p_\theta(x_{t-1}|x_t)) \rightarrow \mathcal{L}_{simple}(\theta) = \|\epsilon_\theta(x_t) - \epsilon_t\|_2^2$
- Classifier-free guidance (CFG): $\hat{\epsilon}(x_t, t, c) = \epsilon_\theta(x_t, t, \emptyset) + s(\epsilon_\theta(x_t, t, c) - \epsilon_\theta(x_t, t, \emptyset))$



Diffusion based Transformer

Large Language Diffusion Models (Nie et al., 2025, Neurips 2025 Oral)

- Forward process: $q_{t|0}(x_t|x_0) = \prod_{i=1}^L q_{t|0}(x_t^i|x_0^i)$
- $q_{t|0}(x_t^i|x_0^i) = \begin{cases} 1 - t, x_t^i = x_0^i, \\ t, x_t^i = M, M \text{ denotes the mask token} \end{cases}$
- Loss (training) : $\mathcal{L}(\theta) = -\mathbb{E}_{t,x_0,x_t} \left[\frac{1}{t} \sum_{i=1}^L \mathbb{1}[x_t^i = M] \log p_{\theta}(x_0^i|x_t) \right]$

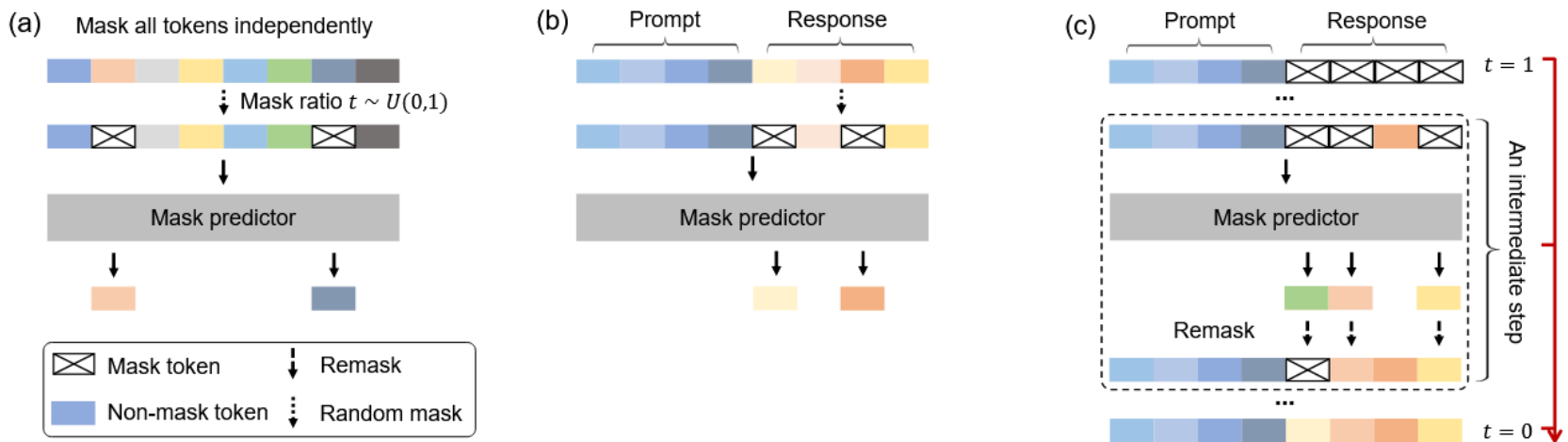
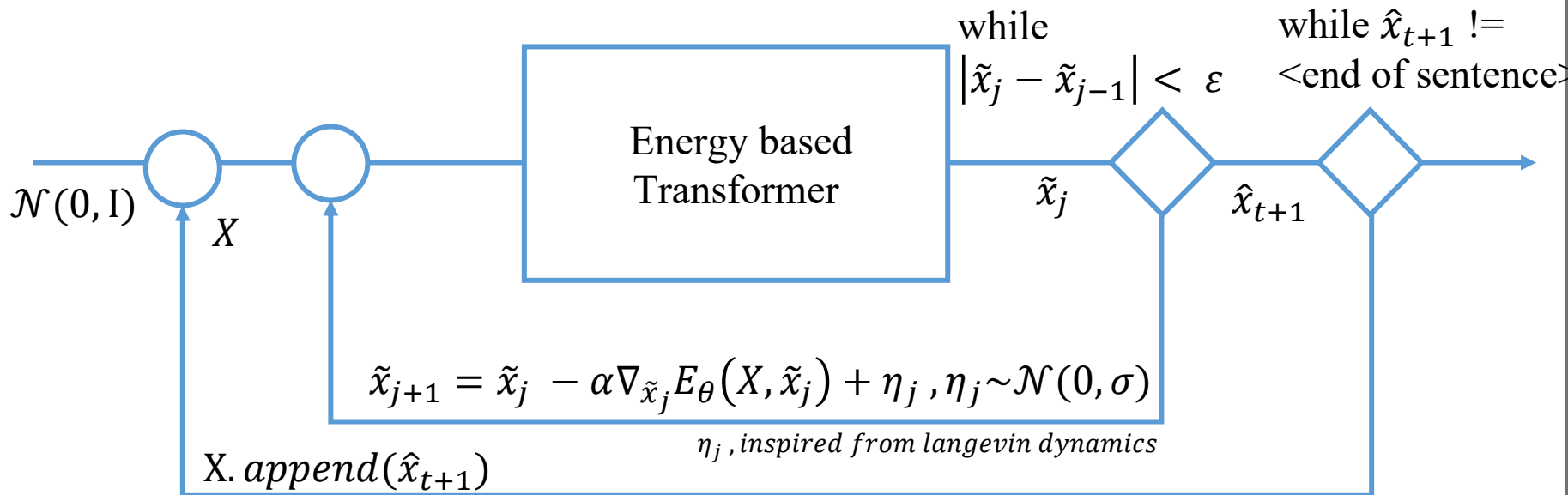


Figure 2. **A Conceptual Overview of LLaDA.** (a) Pre-training. LLaDA is trained on text with random masks applied independently to all tokens at the same ratio $t \sim U[0, 1]$. (b) SFT. Only response tokens are possibly masked. (c) Sampling. LLaDA simulates a diffusion process from $t = 1$ (fully masked) to $t = 0$ (unmasked), predicting all masks simultaneously at each step with flexible remask strategies.

Energy-Based Transformer

- Energy-Based Transformers (Gladstone, Alex et al, 2026) perform energy minimization in a locally **convex representation space** and **automatically stop upon convergence**, enabling scalable unsupervised System 2 Thinking



Energy-Based Transformer

- ENERGY-BASED TRANSFORMERS ARE SCALABLE LEARNERS AND THINKERS (Gladstone, Alexi et al. , ICLR 2026 Oral)

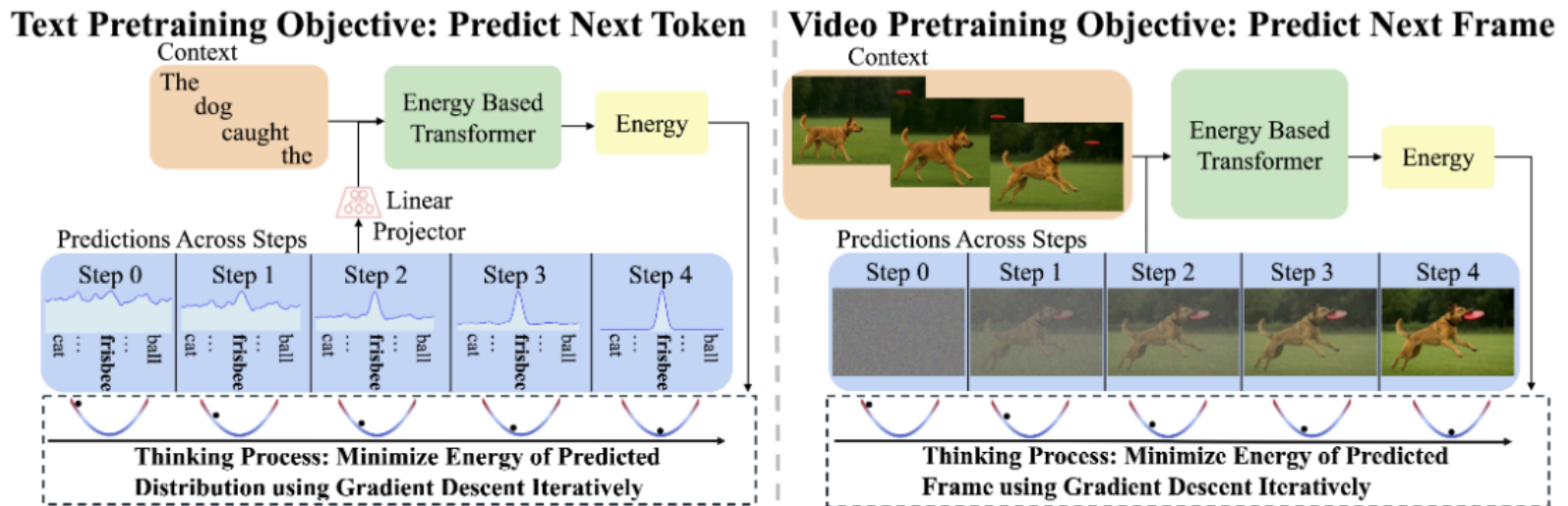
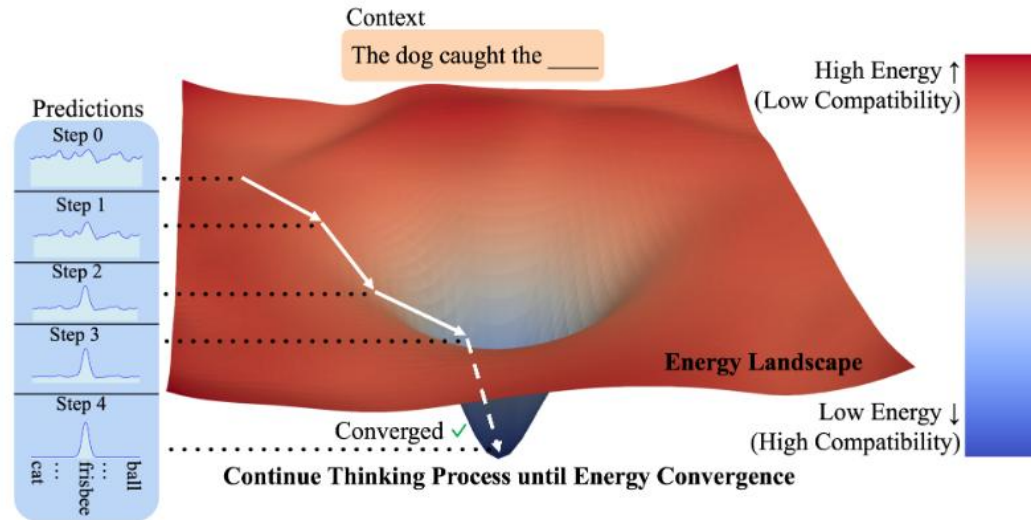


Figure 2: **EBT for Autoregressive Modeling.** Each blue box corresponds to a different prediction at each step of the thinking process, where the initial prediction starts as random. At each step, a new prediction is fed into the model, which gives an energy scalar for the prediction's current *compatibility* (unnormalized likelihood) with the context (Facet 2). Then, the gradient of this energy with respect to the prediction is calculated and used to update the prediction. This gradient descent update is done iteratively to refine the prediction until convergence of the predicted energy, which allows for dynamic use of computation (Facet 1).

Energy-Based Transformer

- ENERGY-BASED TRANSFORMERS ARE SCALABLE LEARNERS AND THINKERS (Gladstone, Alexi et al. , ICLR 2026 Oral)



Algorithm 1: Training

Inputs: Context x , Target y , EBM $E_\theta(x, \hat{y})$

Hparams: Steps M , Step Size α , Loss $J(\cdot)$

- 1 Sample $\hat{y}_0 \sim \mathcal{N}(0, I)$;
 - 2 **for** $i = 0, \dots, M - 1$ **do**
 - 3 | $\hat{y}_{i+1} \leftarrow \hat{y}_i - \alpha \nabla_{\hat{y}_i} E_\theta(x, \hat{y}_i)$;
 - 4 $\mathcal{L} \leftarrow J(\hat{y}_M, y)$;
 - 5 **return** \mathcal{L} , *update* E_θ ;
-

Algorithm 2: Inference with Verification

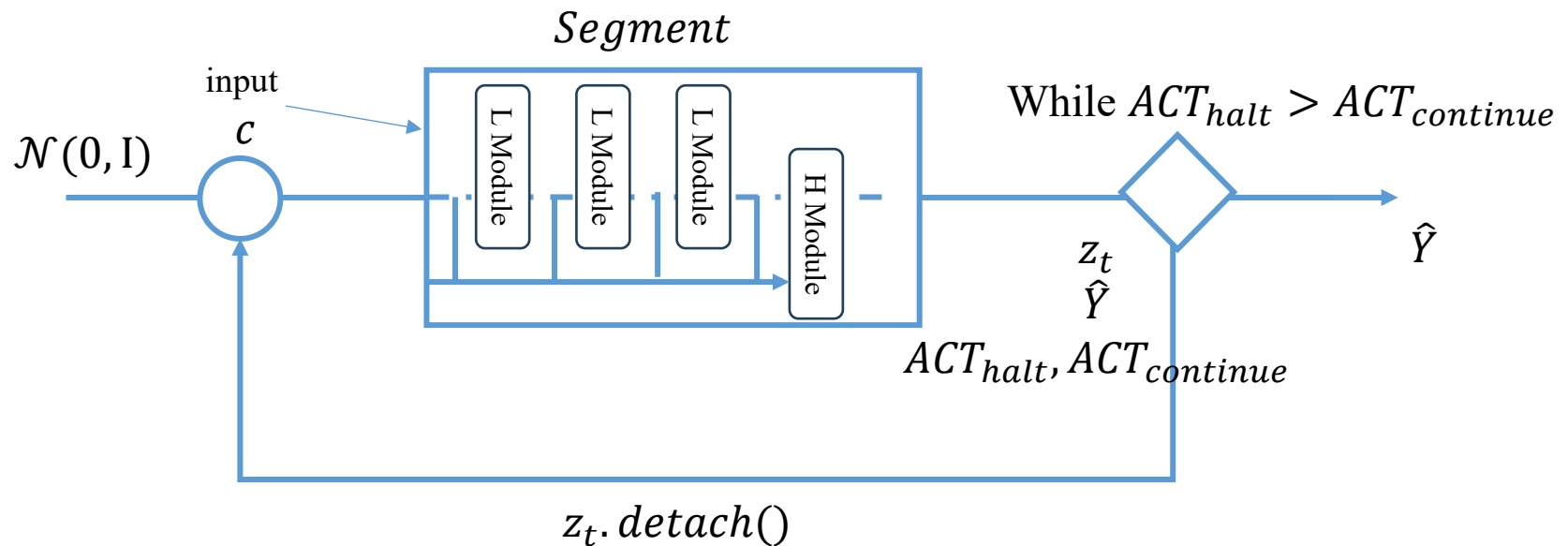
Inputs: Context x , EBM $E_\theta(x, \hat{y})$

Hparams: Steps M , Step Size α , Samples N

- 1 **for** $j = 1, \dots, N$ **do**
 - 2 | Sample $\hat{y}_{0,j} \sim \mathcal{N}(0, I)$;
 - 3 | **for** $i = 0, \dots, M - 1$ **do**
 - 4 | | $\hat{y}_{i+1,j} \leftarrow \hat{y}_{i,j} - \alpha \nabla_{\hat{y}_{i,j}} E_\theta(x, \hat{y}_{i,j})$;
 - 5 **return** $\hat{y}^* = \operatorname{argmin}_j E_\theta(x, \hat{y}_{M,j})$;
-

Hierarchical Reasoning Model

- Inspired by brain dynamics, HRM was designed two Transformer-based modules : **H module** (slower timescales like theta waves) and **L module** (fast timescales like gamma waves)
- **Deep supervision, 1-step Approximate gradient** and **Adaptive Computational Time(ACT)**, the model learns dynamic thinking trajectories end-to-end, without requiring fixed noise schedules(DiT) or predefined convexity assumptions(EBT).



Hierarchical Reasoning Model

■ Hierarchical Reasoning Model (Sapient Intelligence, arXiv 25.06)

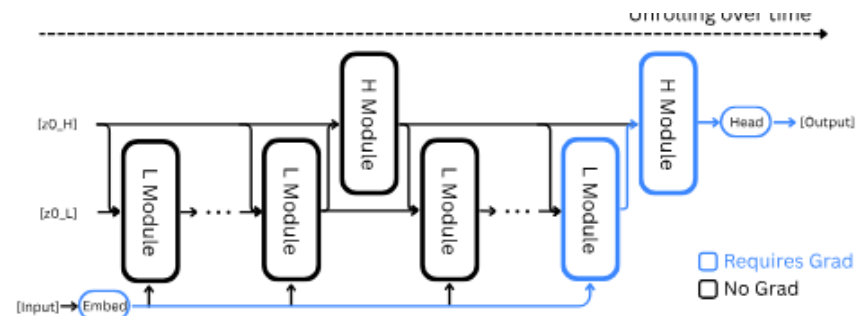
```

def hrm(z, x, n=2, T=2): # hierarchical reasoning
    zH, zL = z
    with torch.no_grad():
        for i in range(nT - 2):
            zL = L_net(zL, zH, x)
            if (i + 1) % T == 0:
                zH = H_net(zH, zL)
    # 1-step grad
    zL = L_net(zL, zH, x)
    zH = H_net(zH, zL)
    return (zH, zL), output_head(zH), Q_head(zH)

def ACT_halt(q, y_hat, y_true):
    target_halt = (y_hat == y_true)
    loss = 0.5*binary_cross_entropy(q[0], target_halt)
    return loss

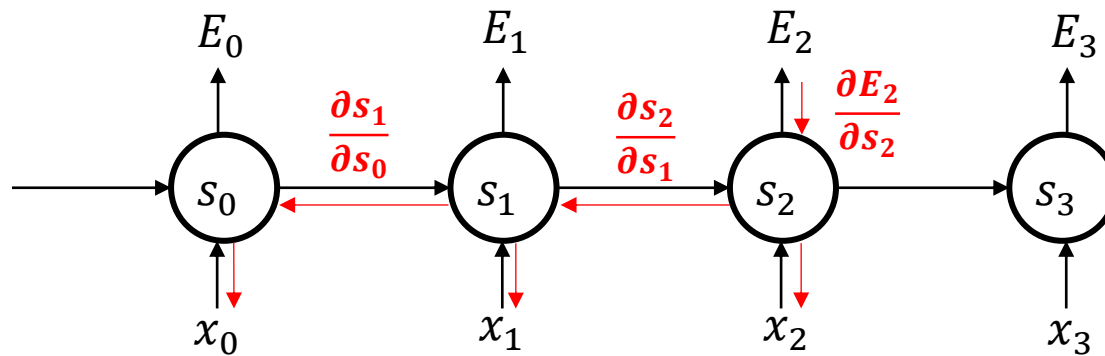
def ACT_continue(q, last_step):
    if last_step:
        target_continue = sigmoid(q[0])
    else:
        target_continue = sigmoid(max(q[0], q[1]))
    loss = 0.5*binary_cross_entropy(q[1], target_continue)
    return loss

# Deep Supervision
for x_input, y_true in train_dataloader:
    z = z_init
    for step in range(N_sup): # deep supervision
        x = input_embedding(x_input)
        z, y_pred, q = hrm(z, x)
        loss = softmax_cross_entropy(y_pred, y_true)
        # Adaptive computational time (ACT) using Q-learning
        loss += ACT_halt(q, y_pred, y_true)
        _, _, q_next = hrm(z, x) # extra forward pass
        loss += ACT_continue(q_next, step == N_sup - 1)
        z = z.detach()
        loss.backward()
        opt.step()
        opt.zero_grad()
        if q[0] > q[1]: # early-stopping
            break
    
```

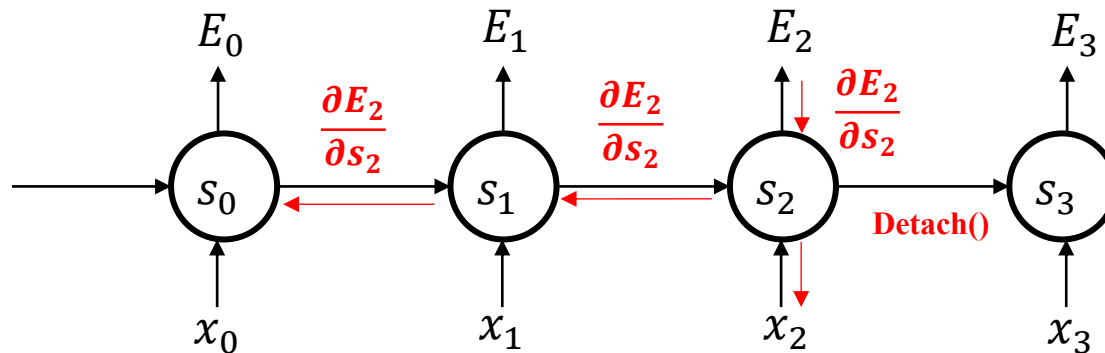


Hierarchical Reasoning Model

- Hierarchical Reasoning Model (Sapient Intelligence, arXiv 25.06)



BPTT(Backpropagation Through Time) (Werbos, 1990)



1-step Approximation (Geng, 2021)

Hierarchical Reasoning Model

■ Deep Equilibrium Model(DEQ): Ideal Segment State of HRM

- The low-level state z_L^* is defined as the **fixed point** of repeated updates:

$$z_L^* = f_L(z_L^*, z_H^{k-1}, \tilde{x}; \theta_L)$$

- Where the iterative application of f_L cause the system to converge to the fixed point z_L^* . When the **H-module** receives this converged result and performs a single update, it is defined as:

$$z_H^k = f_H(z_H^{k-1}, z_L^*, \tilde{x}; \theta_H), \quad z_H^* = F(z_H^{k-1}, \tilde{x}; \theta)$$

- Let $J_F = \frac{\partial F}{\partial z_H}$, At the fixed point, the gradient with respect to parameters satisfies:

$$\frac{\partial z_H^*}{\partial \theta} = J_F \frac{\partial z_H^*}{\partial \theta} + \frac{\partial F}{\partial \theta} \text{ which leads to: } \frac{\partial z_H^*}{\partial \theta} = (I - J_F)^{-1} \frac{\partial F}{\partial \theta}$$

- Using the Neumann series expansion,

$$(I - J_F)^{-1} = 1 + J_F + J_F^2 + J_F^3 \dots, (I - J_F)^{-1} \cong I$$

$$\frac{\partial z_H^*}{\partial \theta_H} \approx \frac{\partial f_H}{\partial \theta_H}, \quad \frac{\partial z_H^*}{\partial \theta_L} \approx \frac{\partial f_H}{\partial z_L^*} \cdot \frac{\partial z_L^*}{\partial \theta_L}, \quad \frac{\partial z_H^*}{\partial \theta_I} \approx \frac{\partial f_H}{\partial z_L^*} \cdot \frac{\partial z_L^*}{\partial \theta_I}. \quad (2)$$

$$\frac{\partial z_L^*}{\partial \theta_L} \approx \frac{\partial f_L}{\partial \theta_L}, \quad \frac{\partial z_L^*}{\partial \theta_I} \approx \frac{\partial f_L}{\partial \theta_I}.$$

Hierarchical Reasoning Model

▪ Deep Supervision

- When the recursive structure of HRM is trained with **deep supervision**:
- **1. Continuous feedback enables iterative updates**, similar to diffusion models that learn from the error at each step.
- **2. Long-term credit assignment (e.g., the vanishing gradient problem)** can be handled approximately and locally, leading to more stable training.

```
# Deep Supervision
for x, y_true in train_dataloader:
    z = z_init
    for step in range(N_supervision):
        z, y_hat = hrm(z, x)

        loss = softmax_cross_entropy(y_hat, y_true)
        z = z.detach()

        loss.backward()
        opt.step()
        opt.zero_grad()
```

Hierarchical Reasoning Model

■ Adaptive computational time (ACT)

- The objective of **ACT** is to enable the model to **autonomously decide how much computation is required for each input**.

• Q-HEAD

- At segment m , the H-module produces Q-values via a Q-head:

$$\widehat{Q} = \sigma(\theta_Q^T z_H^{mNT}) = (\widehat{Q}_{halt}^m, \widehat{Q}_{continue}^m)$$

• Halting Policy

- The halting policy is constrained by a maximum and minimum number of segments, M_{max} **and** M_{min} . With probability ϵ , the minimum number of segments is sampled as $M_{min} \sim \text{Uniform}\{2, \dots, M_{max}\}$, and with probability $1 - \epsilon$, $M_{min} = 1$.

- The model halts when:

- $m > M_{min}$ && $\widehat{Q}_{halt}^m > \widehat{Q}_{continue}^m$ OR $m > M_{max}$

• ACT as Episodic MDP

- State at step m : z^m , Action space : {halt, continue}

- The target values are defined as:

$$\widehat{G}^m = (\widehat{G}_{halt}^m, \widehat{G}_{continue}^m)$$
$$\widehat{G}_{halt}^m = \mathbf{1}\{\hat{y}^m = y\},$$
$$\widehat{G}_{continue}^m = \begin{cases} \widehat{Q}_{halt}^m, & \text{if } m \geq M_{max} \\ \max(\widehat{Q}_{halt}^m, \widehat{Q}_{continue}^m), & \text{otherwise.} \end{cases}$$

• ACT Objective

$$L_{ACT}^m = \text{Loss}(\hat{y}^m, y) + \text{BCE}(\widehat{Q}^m, \widehat{G}^m)$$

Hierarchical Reasoning Model

- Adaptive computational time (ACT)

```
def ACT_halt(q, y_hat, y_true):
    target_halt = (y_hat == y_true)
    loss = 0.5*binary_cross_entropy(q[0], target_halt)
    return loss

def ACT_continue(q, last_step):
    if last_step:
        target_continue = sigmoid(q[0])
    else:
        target_continue = sigmoid(max(q[0], q[1]))
    loss = 0.5*binary_cross_entropy(q[1], target_continue)
    return loss
```

```
# Deep Supervision
for x_input, y_true in train_dataloader:
    z = z_init
    for step in range(N_sup): # deep supervision
        x = input_embedding(x_input)
        z, y_pred, q = hrm(z, x)
        loss = softmax_cross_entropy(y_pred, y_true)
        # Adaptive computational time (ACT) using Q-learning
        loss += ACT_halt(q, y_pred, y_true)
        _, _, q_next = hrm(z, x) # extra forward pass
        loss += ACT_continue(q_next, step == N_sup - 1)
    z = z.detach()
    loss.backward()
    opt.step()
    opt.zero_grad()
    if q[0] > q[1]: # early-stopping
        break
```

Hierarchical Reasoning Model

▪ Less is More: Recursive Reasoning with Tiny Networks (Samsung, arXiv, 2510)

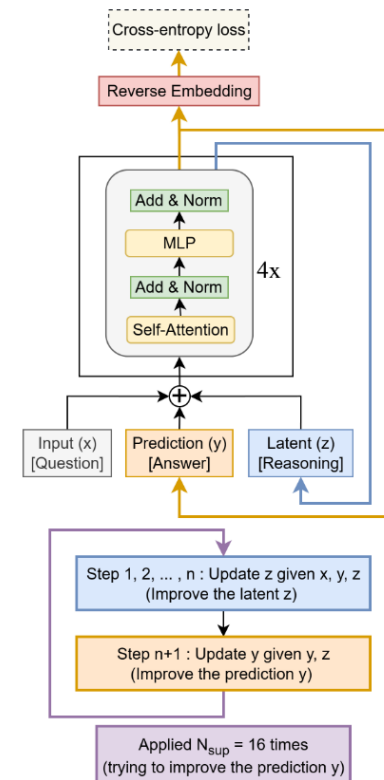
- 1. Limitations of the **fixed-point** assumption, particularly the reliance on one-step approximation.
- 2. The H-module and L-module architecture contributes only **marginally** to performance.
- 3. **Deep supervision** is the primary driver of strong performance.
- 4. Q-learning can be replaced with a **simple binary classifier**.

```
def latent_recursion(x, y, z, n=6):
    for i in range(n): # latent reasoning
        z = net(x, y, z)
        y = net(y, z) # refine output answer
    return y, z

def deep_recursion(x, y, z, n=6, T=3):
    # recursing T-1 times to improve y and z (no gradients needed)
    with torch.no_grad():
        for j in range(T-1):
            y, z = latent_recursion(x, y, z, n)
    # recursing once to improve y and z
    y, z = latent_recursion(x, y, z, n)
    return (y.detach(), z.detach()), output_head(y), Q_head(y)

# Deep Supervision
for x_input, y_true in train_dataloader:
    y, z = y_init, z_init
    for step in range(N_supervision):
        x = input_embedding(x_input)
        (y, z), y_hat, q_hat = deep_recursion(x, y, z)
        loss = softmax_cross_entropy(y_hat, y_true)
        loss += binary_cross_entropy(q_hat, (y_hat == y_true))
        loss.backward()
        opt.step()
        opt.zero_grad()
        if q_hat > 0: # early-stopping
            break
```

Figure 3. Pseudocode of Tiny Recursion Models (TRMs).



Hierarchical Reasoning Model

- Less is More: Recursive Reasoning with Tiny Networks (Samsung, arXiv, 2510)

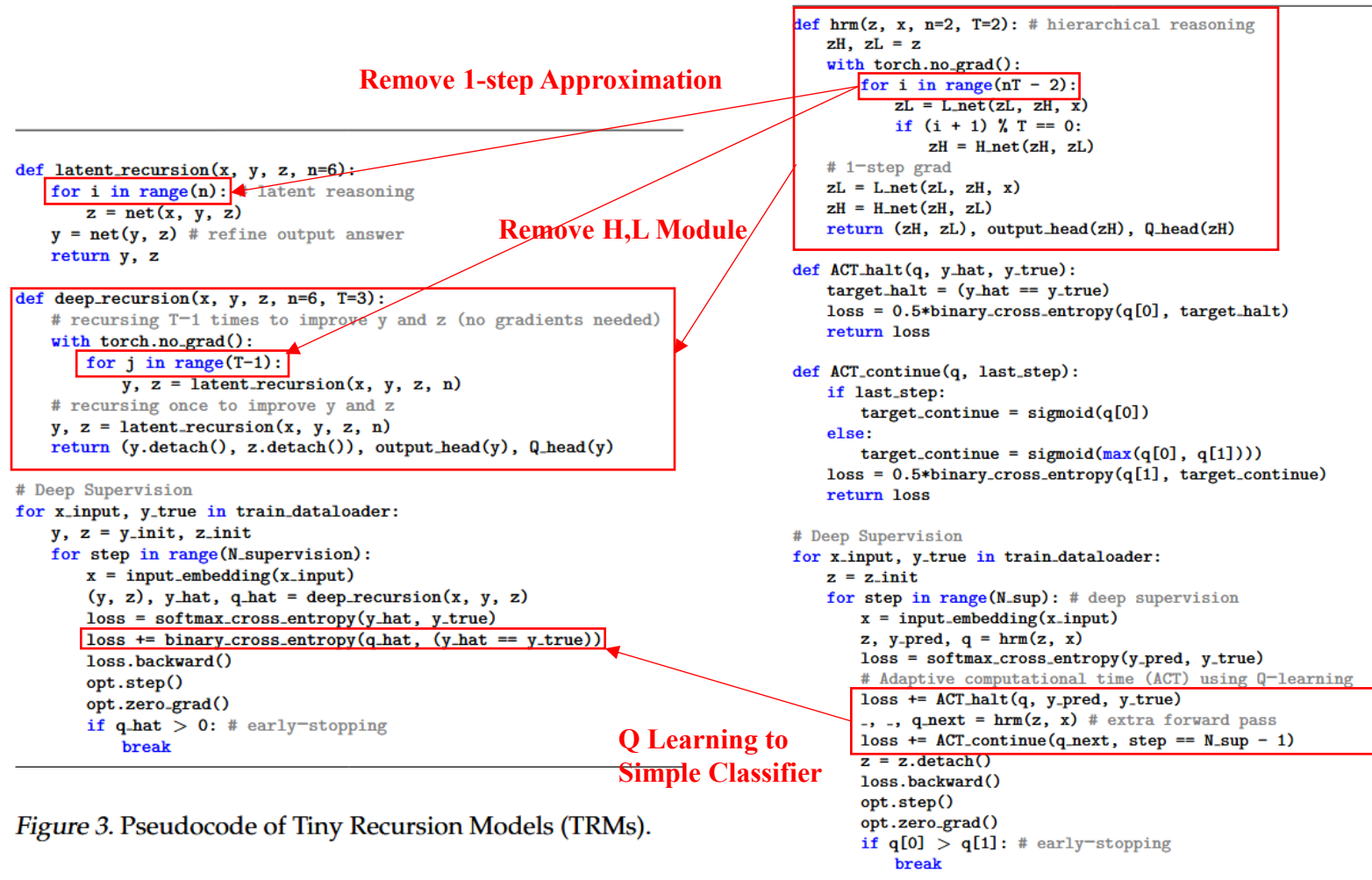
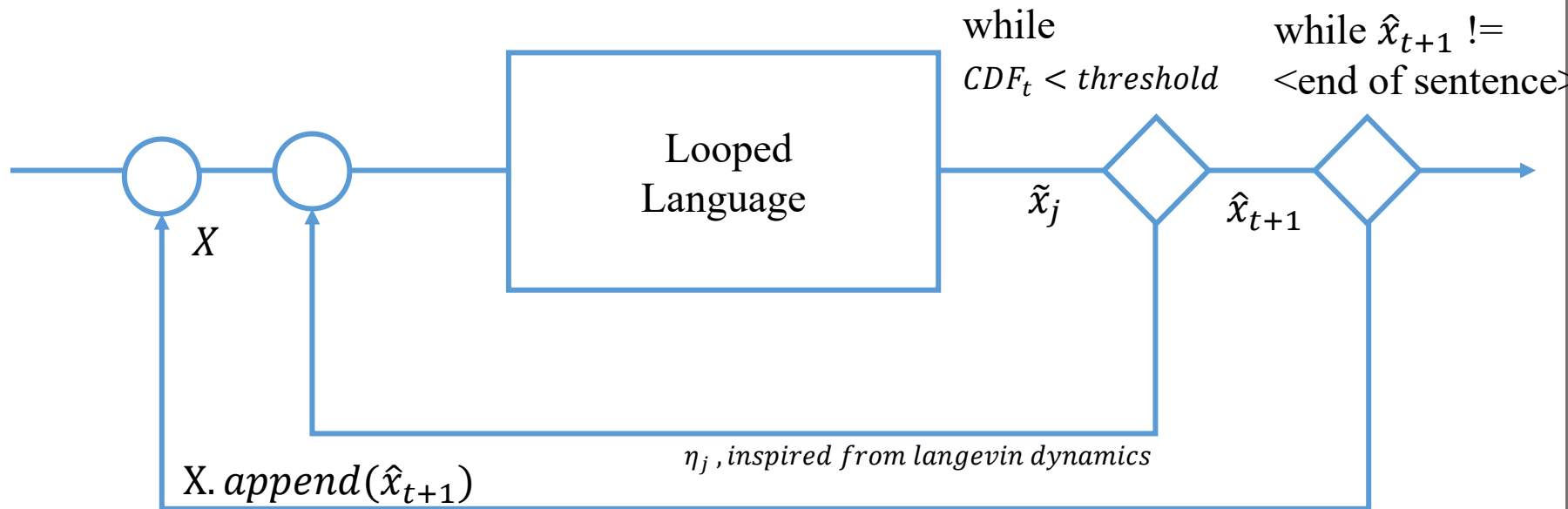


Figure 3. Pseudocode of Tiny Recursion Models (TRMs).

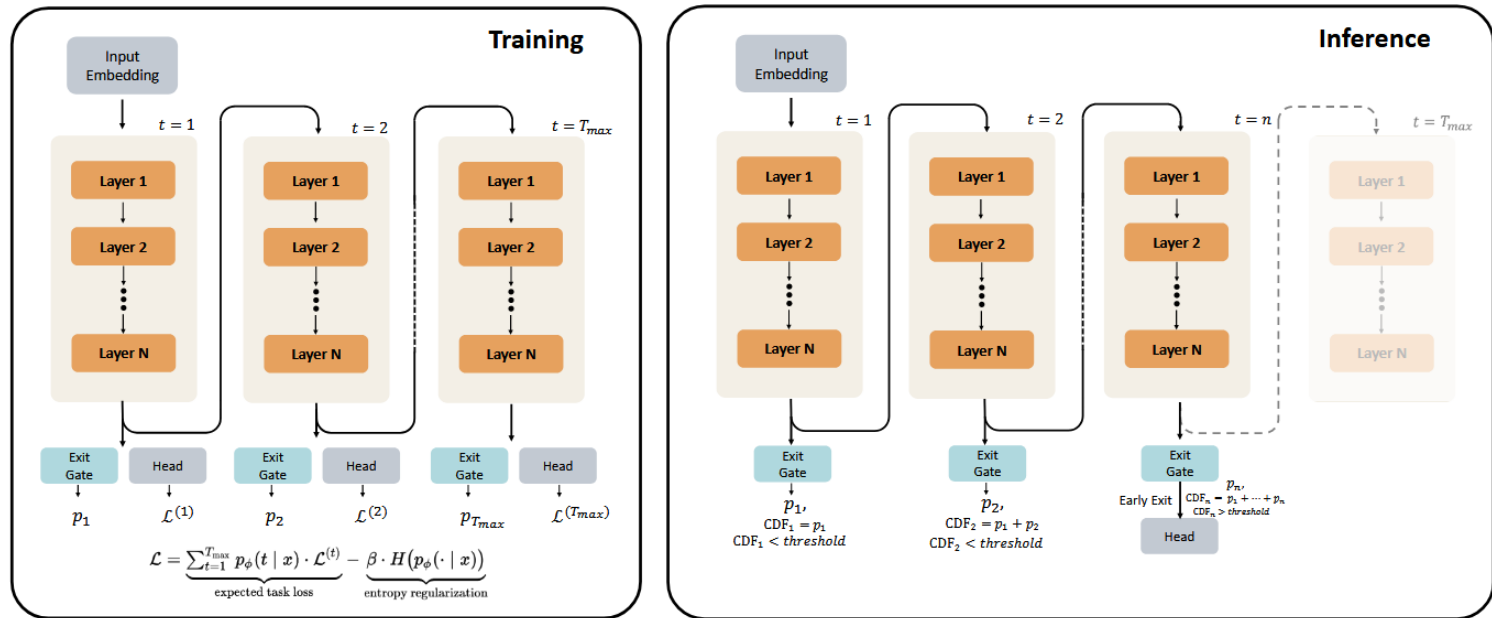
Scaling Latent Reasoning

- Scaling Latent Reasoning via Looped Language Models (ByteDance Seed, Zhu et al, arXiv, 25.10)



Scaling Latent Reasoning

- Scaling Latent Reasoning via Looped Language Models (ByteDance Seed, Zhu et al, arXiv, 25.10)



$$F(\cdot) := \text{lmhead} \circ \mathcal{M}^L \circ \text{emb}(\cdot), \quad \mathcal{M}^L(\cdot) := \mathcal{T}_{\theta_L} \circ \dots \circ \mathcal{T}_{\theta_1}(\cdot)$$

$$F^{(t)}(\cdot) = \text{lmhead} \circ \underbrace{\mathcal{M}^L \circ \mathcal{M}^L \circ \dots \circ \mathcal{M}^L}_{t \text{ iterations}} \circ \text{emb}(\cdot). \quad (1)$$

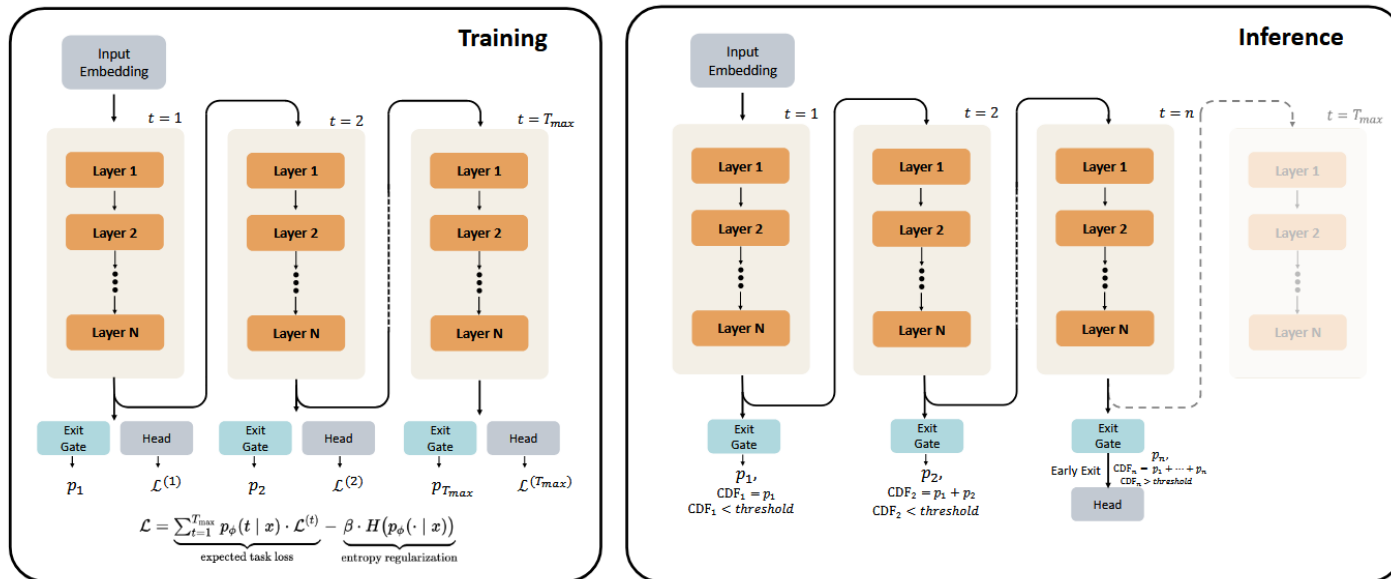
Scaling Latent Reasoning

Scaling Latent Reasoning via Looped Language Models (ByteDance Seed, Zhu et al, arXiv, 25.10)

- Exit probability: $\lambda_t(x) = \sigma(\text{Linear}_\phi(h^{(t)})) \in (0,1)$
- Not exiting in the first t steps : $S_t(x) = \prod_{j=1}^t (1 - \lambda_j(x)), S_0(x) = 1$
- $\tilde{p}_t(x) = \lambda_t(x)S_{t-1}(x), t = 1, \dots, T_{max} - 1$

$$p_\phi(t | x) = \begin{cases} \tilde{p}_t(x), & t = 1, \dots, T_{max} - 1, \\ S_{T_{max}-1}(x), & t = T_{max}, \end{cases} \quad \text{so that} \quad \sum_{t=1}^{T_{max}} p_\phi(t | x) = 1. \quad (3)$$

$$\text{CDF}(n | x) = \sum_{t=1}^n p_\phi(t | x) = 1 - \prod_{j=1}^n (1 - \lambda_j(x)), \quad n < T_{max}, \quad \text{CDF}(T_{max} | x) = 1.$$



Scaling Latent Reasoning

- **Scaling Latent Reasoning via Looped Language Models** (ByteDance Seed, Zhu et al, arXiv, 25.10)
 - **Stable Training** (Stage 1a, Stage 1b) : Joint training of the LM and the gating mechanism for stable adaptive computation.
 - **CT Annealing** (Stage 2) : Freeze the LM and train only the gate to refine depth allocation.
 - **LongCT** (Stage 3) : Extend the context window using long-sequence training.
 - **Mid-training** (Stage 4) : Improve overall model capability using high-quality curated datasets.
 - **Reasoning SFT** (Post Training): Enhance reasoning performance through supervised fine-tuning on reasoning-focused data.

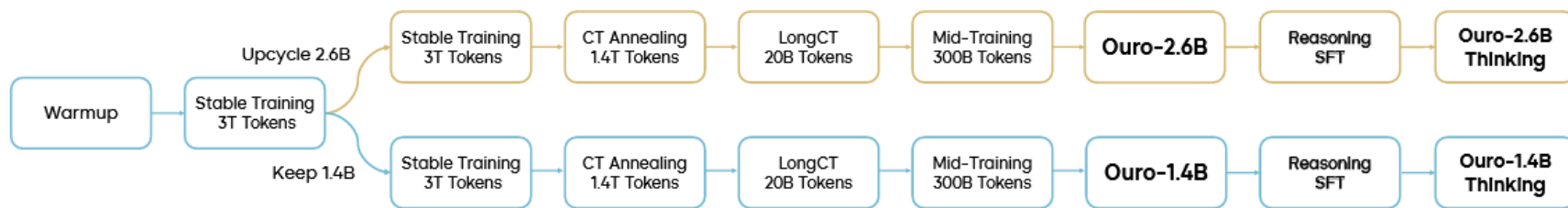
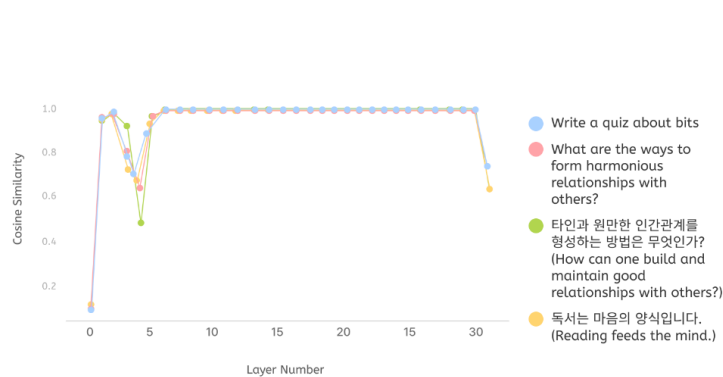


Figure 4 End-to-end Ouro training pipeline: shared warmup → Stable Training → forks into a 1.4B retained path and a 2.6B upcycled path → four shared stages → Reasoning SFT to produce Ouro-Thinking.

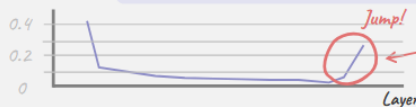
Deep Up-Scaling

From an analysis of the roles of Transformer blocks in LLMs, the early blocks perform representation in a latent space (**Enrichment Phase**), the middle blocks perform thinking (**Plateau Phase**), and the later blocks change the representation back to the original space (**Extraction Phase**) (Tikhonov et al., CIKM 2025).



Observation: Final Layer Hidden States Jump

$$\text{Displacement of hidden state } \Psi_\ell = 1 - \frac{1}{2}(1 + S_C(\mathbf{h}_{\ell-1}, \mathbf{h}_\ell))$$



Assumption: Undesirable property of Transformer-LMs

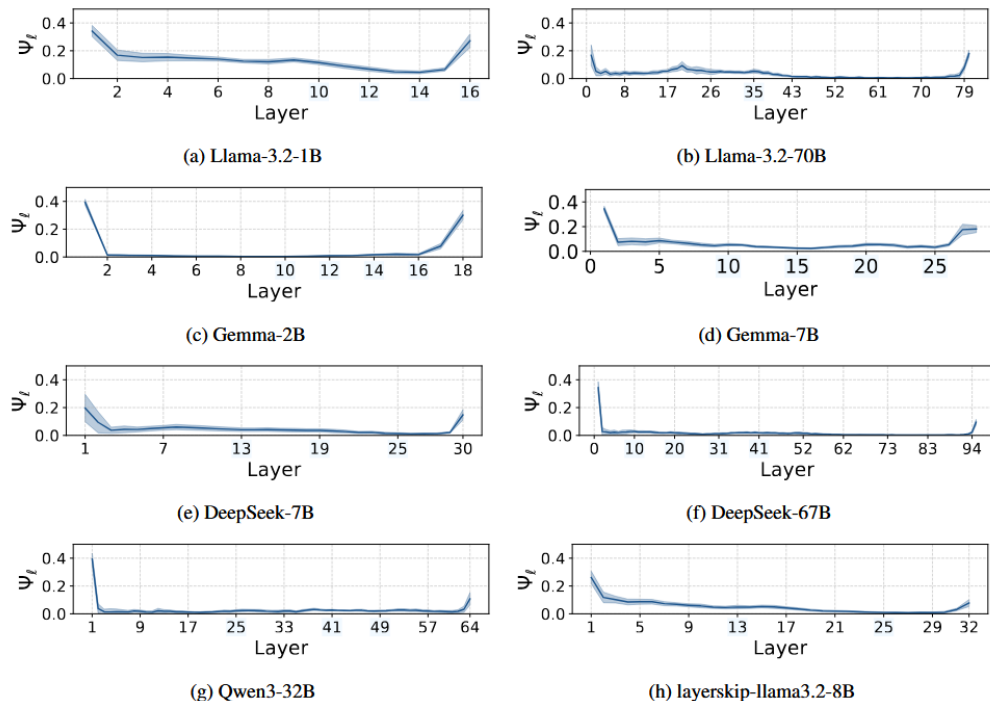
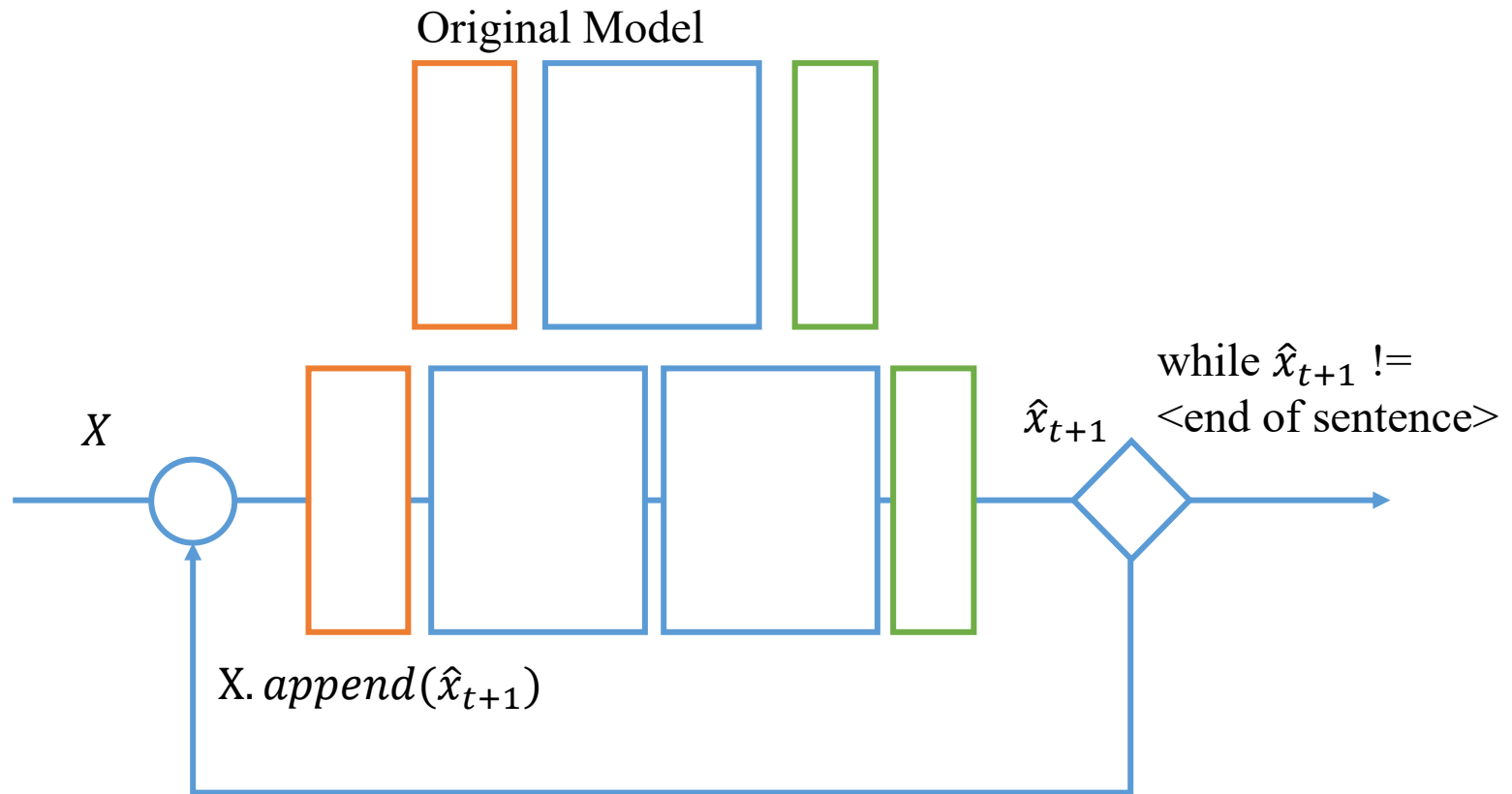


Figure 2: Layer-wise hidden state displacement Ψ_ℓ for next-word prediction on 100 samples from the LAMBADA dataset. Across all model architectures, the displacement at the final layer tends to be larger than that of the middle layers.

Deep Up-Scaling

- **Upstage** (Upstage, NAACL 2024) achieved performance improvements by applying **Deep Up-Scaling** to the LLaMA 2.0 model. Similarly, **Mi:dm 2.0** (KT, arXiv 2601) applies **Deep Up-Scaling** to its internally pretrained model, which is submitted to 독파모.



Deep Up-Scaling

▪ Deep Up-Scaling

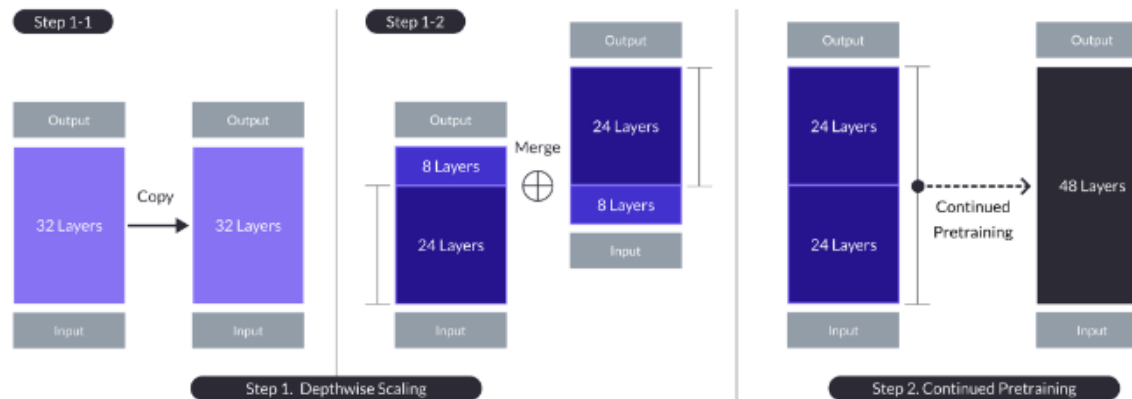
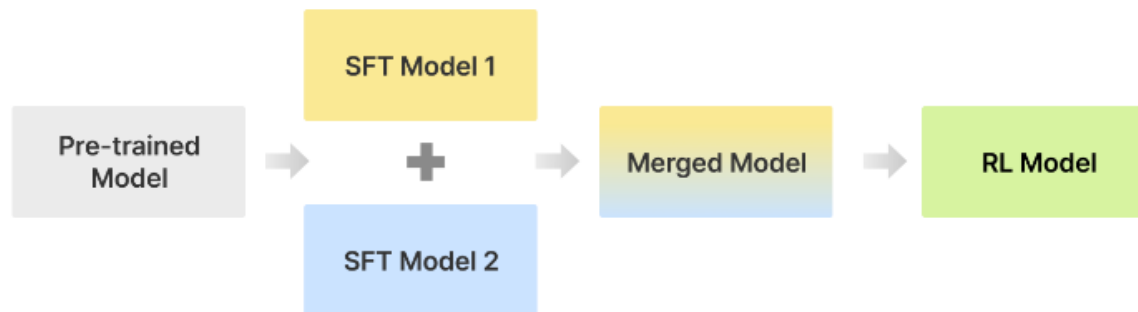
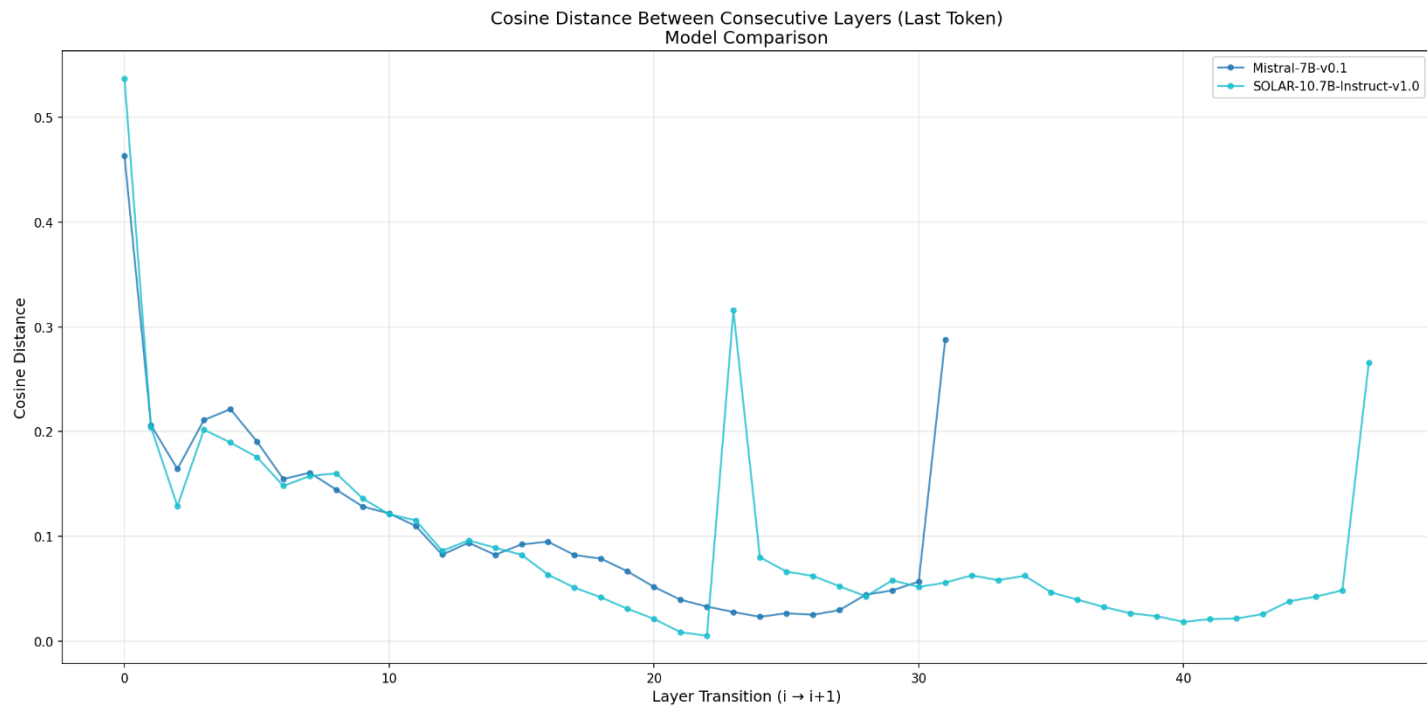


Figure 1: Depth up-scaling for the case with $n = 32$, $s = 48$, and $m = 8$. Depth up-scaling is achieved through a dual-stage process of depthwise scaling followed by continued pretraining.



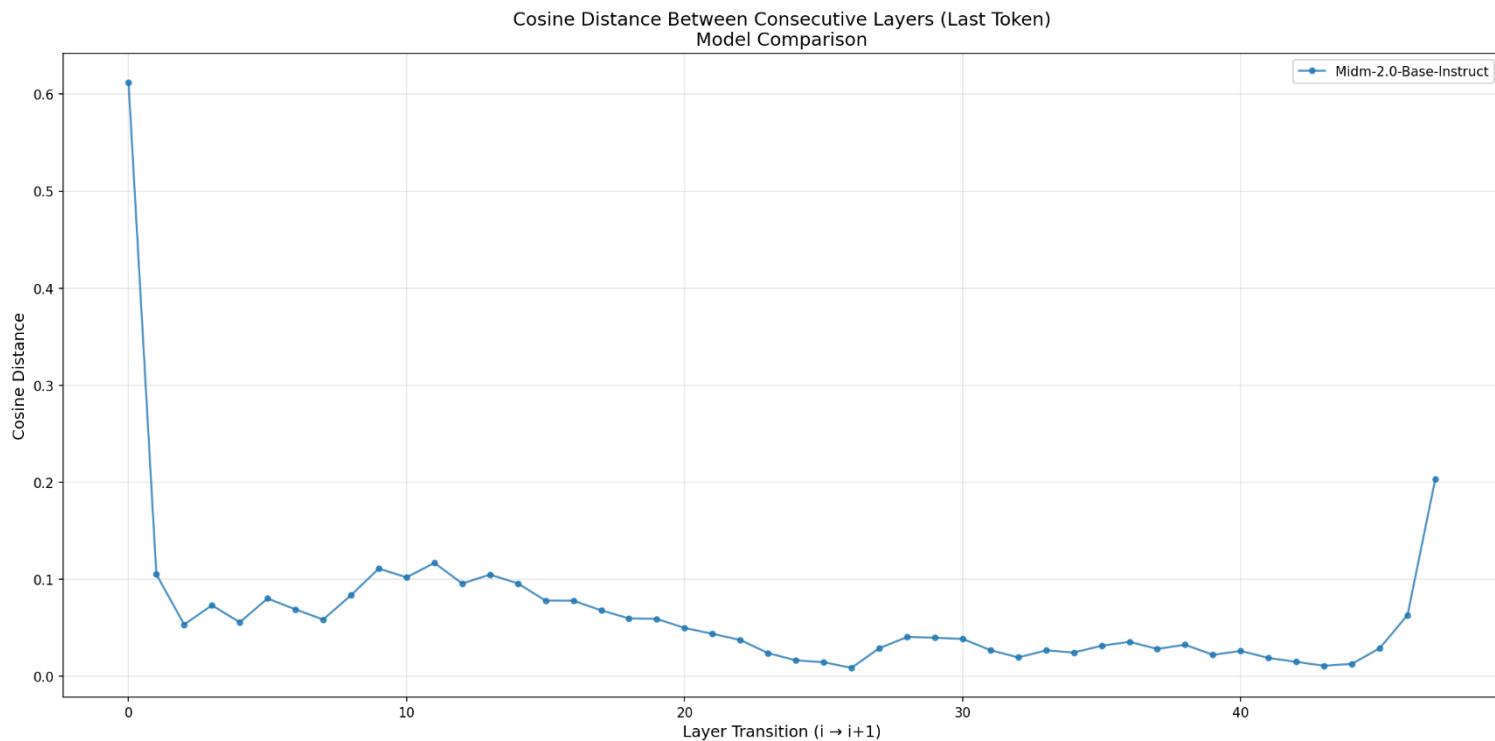
Deep Up-Scaling

■ Deep Up-Scaling



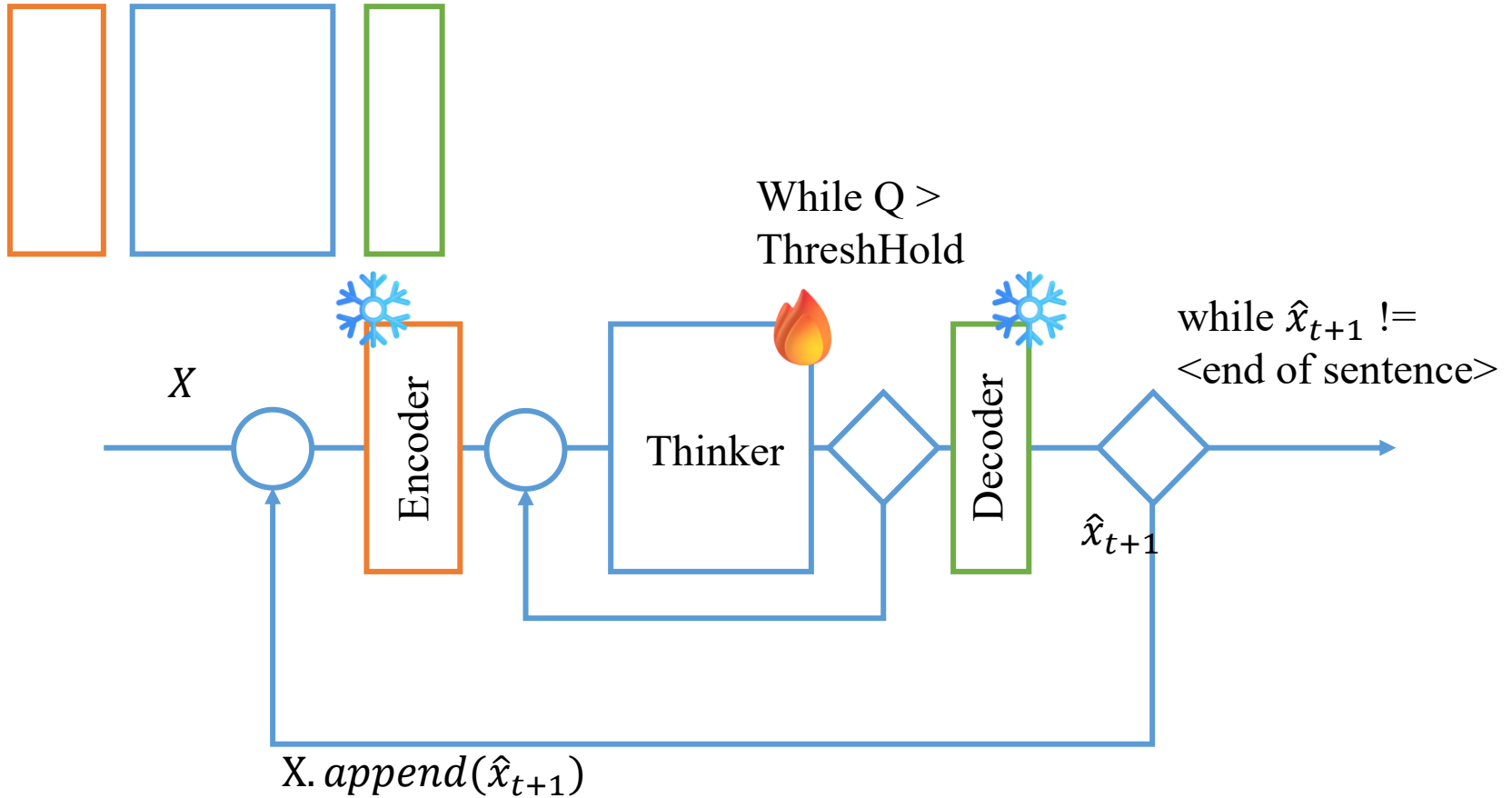
Deep Up-Scaling

■ Deep Up-Scaling

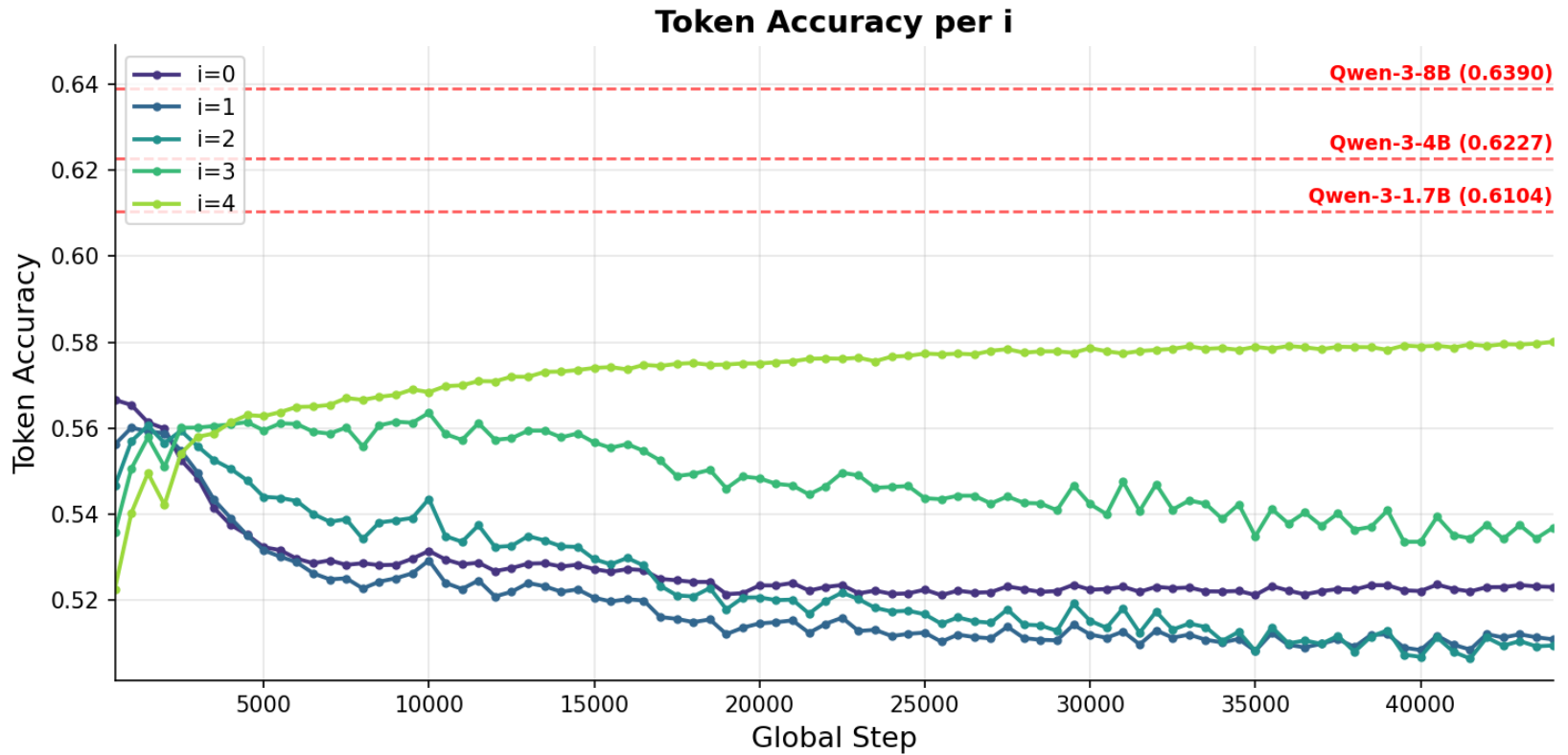


Future Work

Original Model



Future Work



References

- [1] He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
- [2] Liu, Jian-wei, Bing-rong Xu, and Zhi-yan Song. "A Survey of Recursive and Recurrent Neural Networks." *arXiv preprint arXiv:2510.17867* (2025).
- [3] Kahneman, Daniel. *Thinking, fast and slow*. macmillan, 2011.
- [4] Snell, Charlie Victor, et al. "Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning." *The Thirteenth International Conference on Learning Representations*. 2025.
- [5] Hoffmann, Jordan, et al. "Training compute-optimal large language models." *arXiv preprint arXiv:2203.15556* (2022).
- [6] Leviathan, Yaniv, Matan Kalman, and Yossi Matias. "Fast inference from transformers via speculative decoding." *International Conference on Machine Learning*. PMLR, 2023.
- [7] Wang, Guan, et al. "Hierarchical Reasoning Model." *arXiv preprint arXiv:2506.21734* (2025).
- [8] Gladstone, Alexi, et al. "Energy-Based Transformers are Scalable Learners and Thinkers." *arXiv preprint arXiv:2507.02092* (2025).
- [9] Shao, Zhihong, et al. "Deepseekmath: Pushing the limits of mathematical reasoning in open language models." *arXiv preprint arXiv:2402.03300* (2024).
- [10] Peebles, William, and Saining Xie. "Scalable diffusion models with transformers." *Proceedings of the IEEE/CVF international conference on computer vision*. 2023.
- [11] Nie, Shen, et al. "Large Language Diffusion Models." Proceedings of the Thirty-Ninth Annual Conference on Neural Information Processing Systems, 2025
- [12] Gladstone, Alexi, et al. *Energy-Based Transformers Are Scalable Learners and Thinkers*. Proceedings of the International Conference on Learning Representations (ICLR 2026), 2026. Oral Presentation.
- [13] Geng, Zhengyang, Xin-Yu Zhang, Shaojie Bai, Yisen Wang, and Zhouchen Lin. "On Training Implicit Models." *Advances in Neural Information Processing Systems*, 2021.
- [14] **Werbos, Paul J.** "Backpropagation Through Time: What It Does and How to Do It." *Proceedings of the IEEE*, vol. 78, no. 10, 1990, pp. 1550–1560. **IEEE**, <https://doi.org/10.1109/5.58337>.
- [15] Jolicoeur-Martineau, Alexia.(Samsung) "Less is more: Recursive reasoning with tiny networks." *arXiv preprint arXiv:2510.04871* (2025).
- [16] Tikhonov, Pavel, and Dmitry Ilvovsky. "Frozen in the Middle: Hidden States Remain Unchanged Across Intermediate Layers of Language Models." *Proceedings of the 34th ACM International Conference on Information and Knowledge Management*. 2025.
- [17] Kim, Sanghoon, et al(Upstage). "**SOLAR 10.7B: Scaling Large Language Models with Simple yet Effective Depth Up-Scaling.**" *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, vol. 6,
- [18] Shin, Donghoon(KT), et al. "Mi: dm 2.0 Korea-centric Bilingual Language Models." *arXiv preprint arXiv:2601.09066* (2026).
- [19] Geiping, Jonas, et al. "Scaling up test-time compute with latent reasoning: A recurrent depth approach." *arXiv preprint arXiv:2502.05171* (2025).

Q&A