



上海交通大学  
SHANGHAI JIAO TONG UNIVERSITY



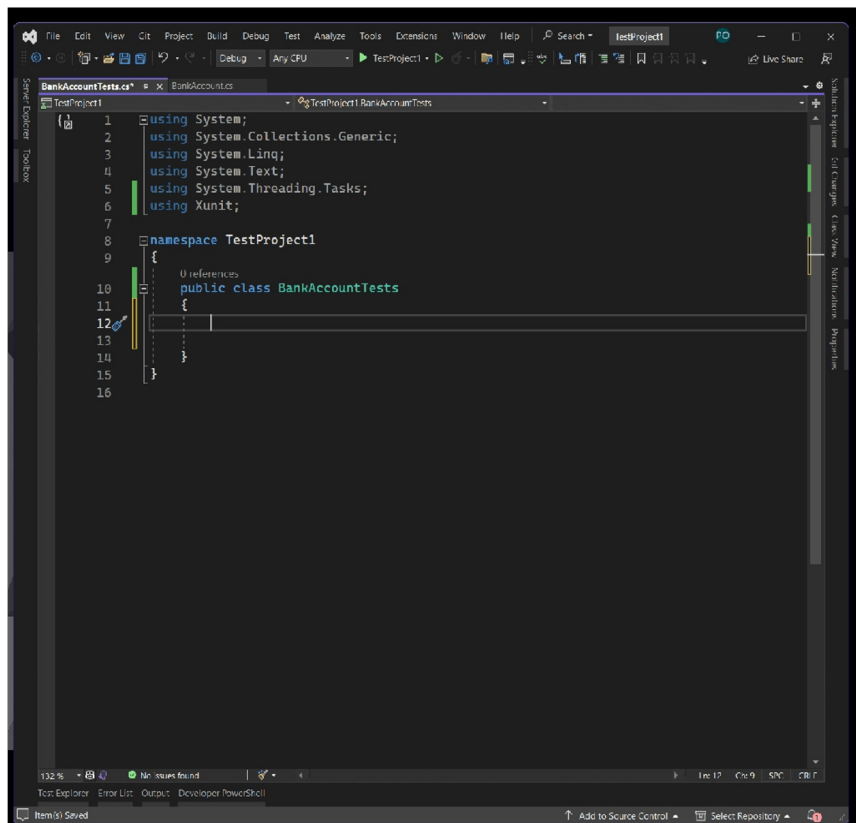
清源研究院  
QING YUAN RESEARCH INSTITUTE

# AIGC大模型的高效推理方法

邓志杰

上海交通大学 清源研究院

# 背景：大语言模型成为广泛下游任务的基石



Code generation



Agents

User: What is funny about this image? Describe it panel by panel.

Source: hmmm (Reddit)

GPT-4: The image shows a package for a "Lightning Cable" adapter with three panels.

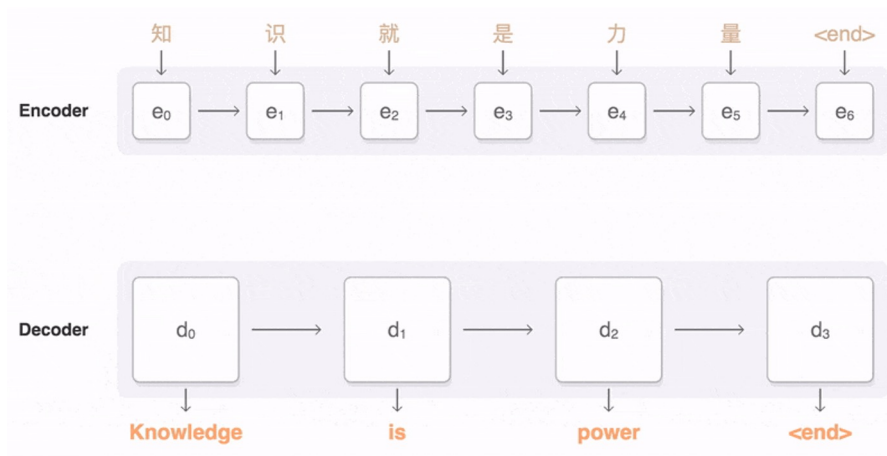
Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.

Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.

Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.

The humor in this image comes from the absurdity of plugging a large, outdated VGA connector into a small, modern smartphone charging port.

(Multimodal) Chatbot



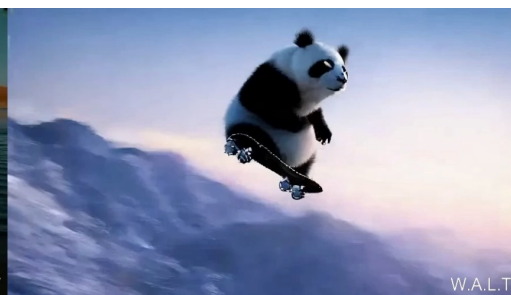
Machine translation

# 背景：大扩散模型正在形成“世界模拟器”



*Sora by OpenAI*

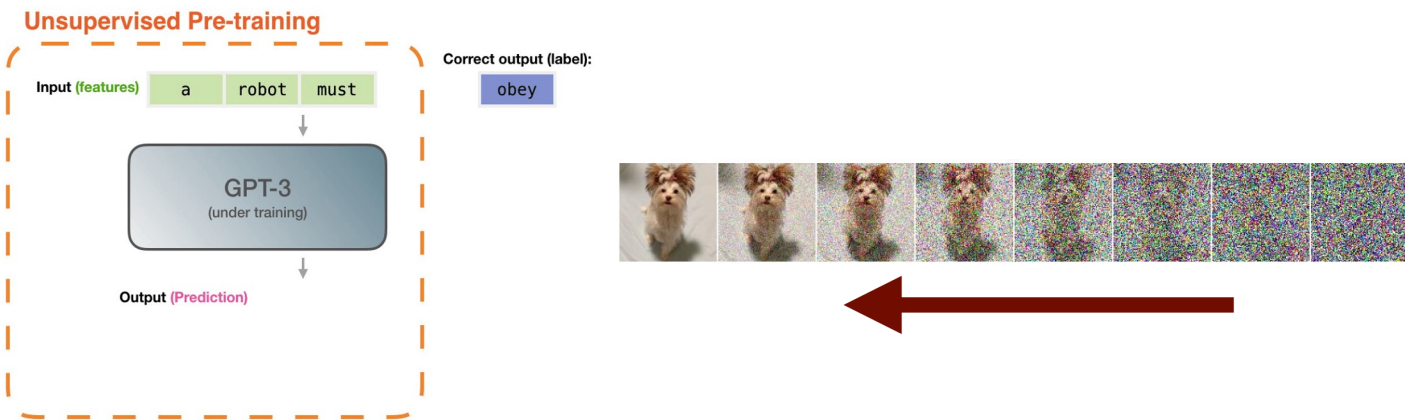
*Vidu by ShengShu*



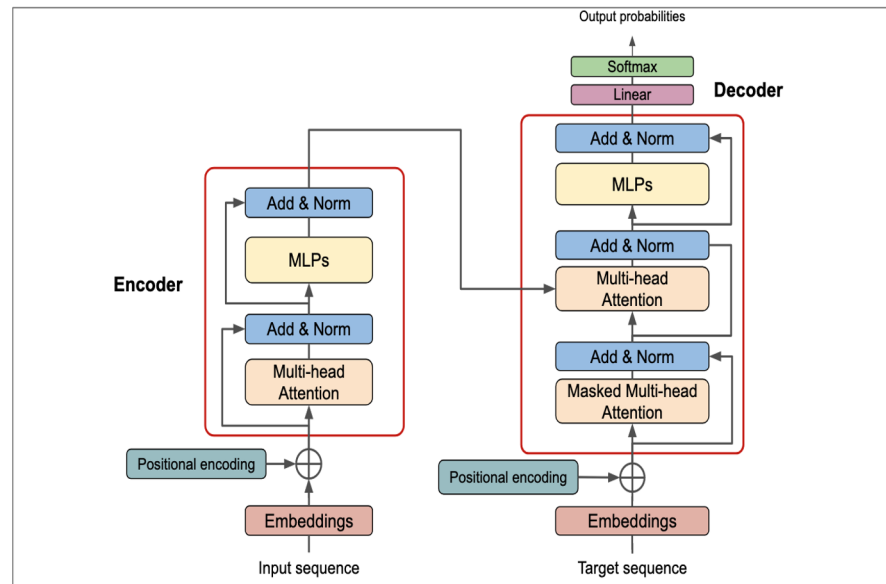
*W.A.L.T*



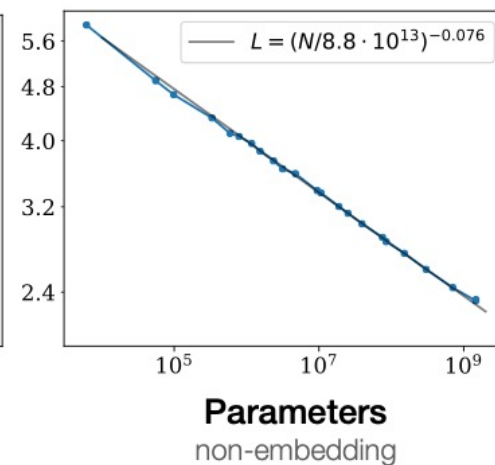
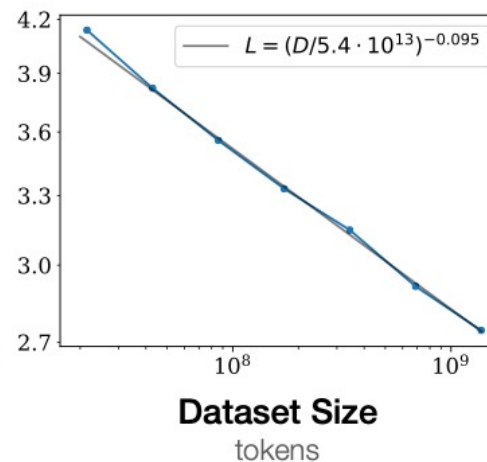
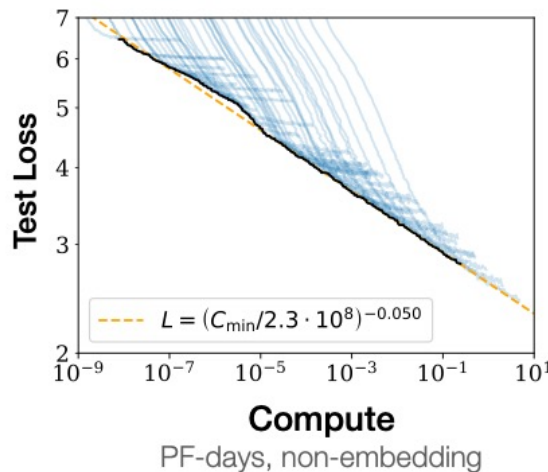
# AIGC大模型的趋势



学习方式存在两个主流



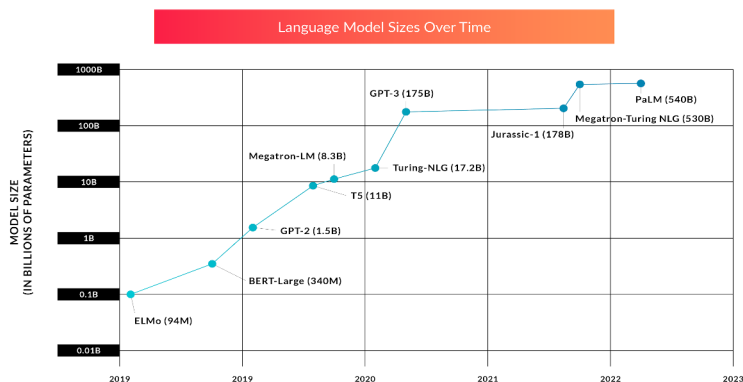
架构趋于统一



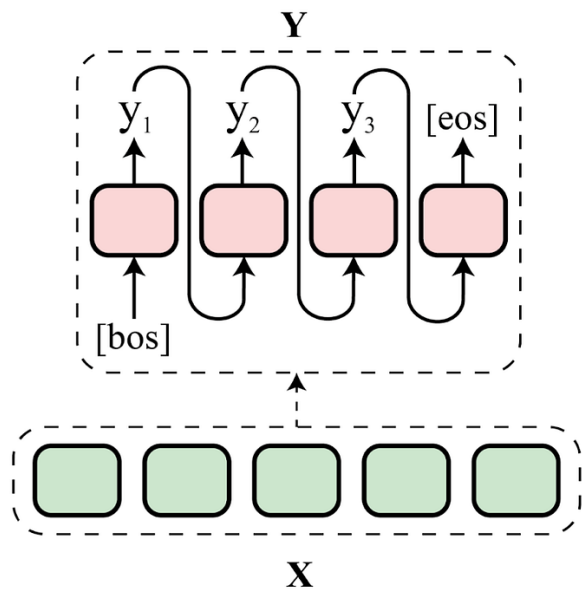
Scaling law持续发挥作用



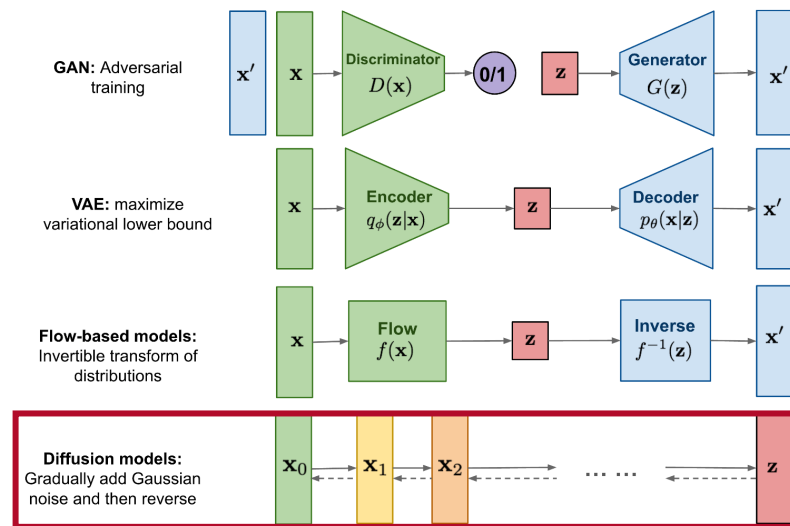
# AIGC大模型推理低效，导致高部署成本、差用户体验



模型本身的大尺寸（更多的 flops、内存占用）



每个新token都需要一次模型前传



每个去噪步都需要一次模型前传

低效的顺序推理过程



# AIGC大模型的高效推理方法

- 大语言模型的并行推理方法（15min）
  - 投机解码
  - 一致性蒸馏
- 大扩散模型的低步数推理方法（15min）
  - 采样器设计
  - 一致性蒸馏
- 大模型架构、序列状态、缓存等方面的优化方法（5min）

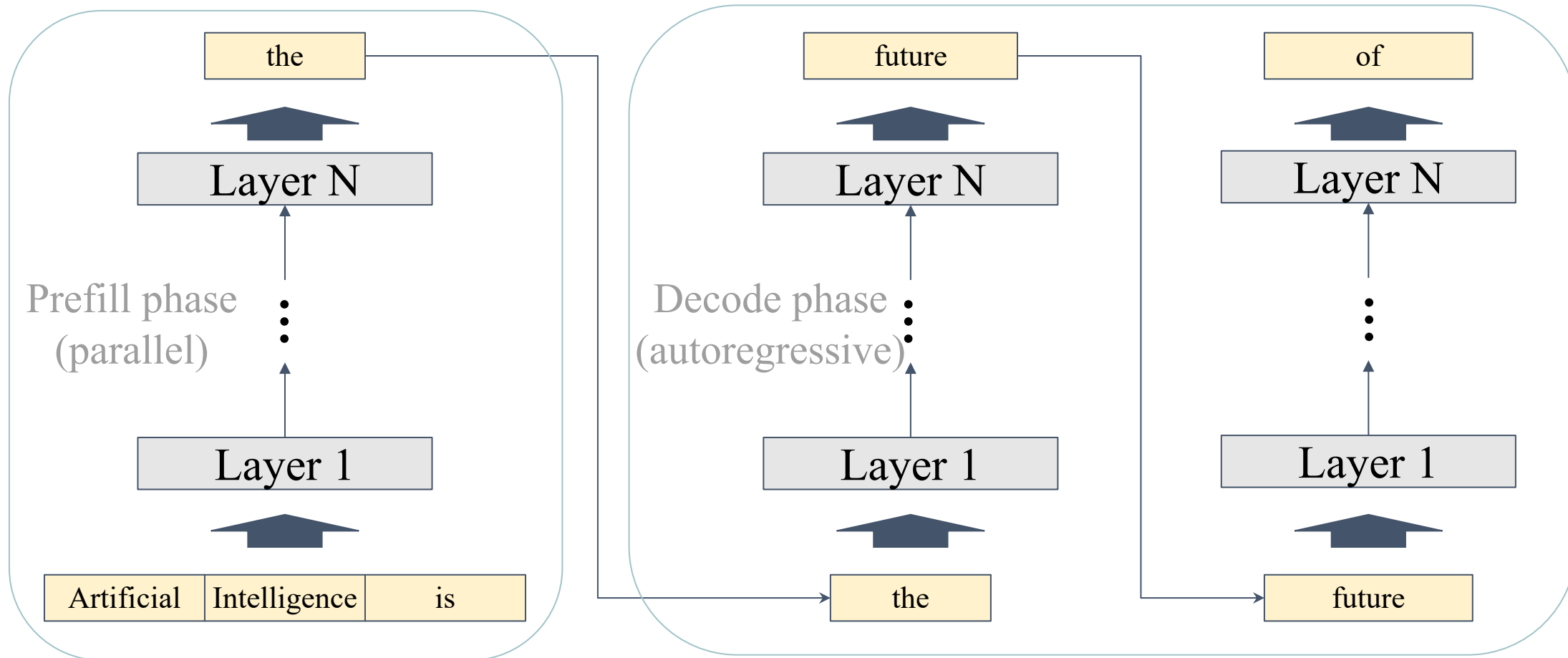


# AIGC大模型的高效推理方法

- 大语言模型的并行推理方法（15min）
  - 投机解码
  - 一致性蒸馏
- 大扩散模型的低步数推理方法（15min）
  - 采样器设计
  - 一致性蒸馏
- 大模型架构、序列状态、缓存等方面的优化方法（5min）



# 大语言模型的推理过程



重复生成过程直达到达到最大生成长度(e.g., 2048 tokens) 或生成出结束符(e.g., “<|end of sequence|>”)

# KV Cache

Output

the

Layer N

Artificial	-0.2	0.1	-1.1
Intelligence	0.9	0.7	0.2
is	-0.1	-0.3	0.1

⋮

Layer 1

Artificial	-0.1	0.3	1.2
Intelligence	0.7	-0.4	0.8
is	0.2	-0.1	1.1



Input

Artificial Intelligence is

future

Layer N

the 

-1.1	0.5	0.4
------	-----	-----

⋮

Layer 1

the 

-0.7	0.1	-0.2
------	-----	------



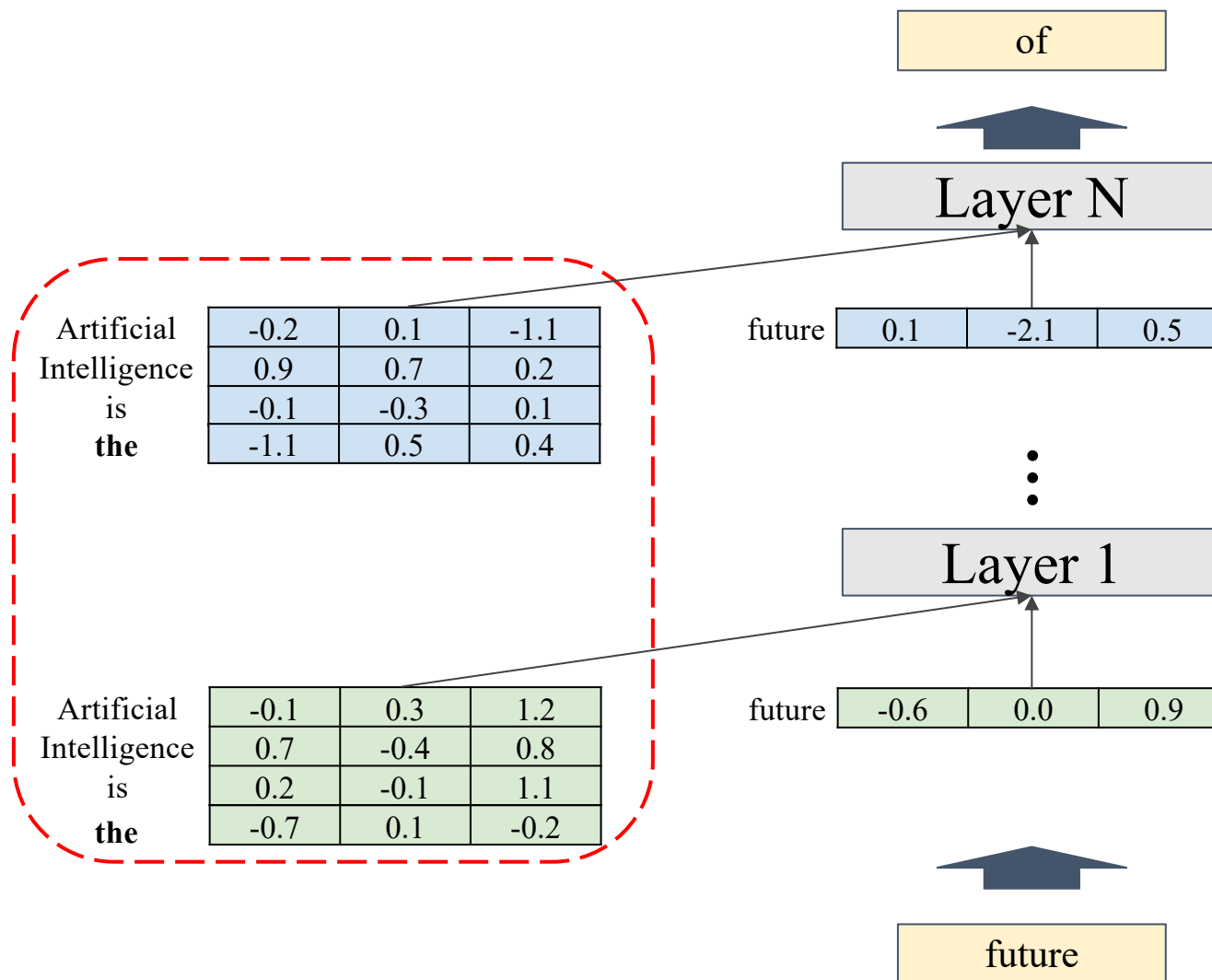
the

# KV Cache

Output

KV Cache

Input





# 大语言模型推理过程的两个特点

- **Prefill phase**计算密集，**decode phase**是**memory-IO bound**
  - between High Bandwidth Memory (HBM) and Static Random Access Memory (SRAM)
- 得益于高效的并行计算，**prefill**几个**tokens**的时间与后续每个**token**生成时间接近

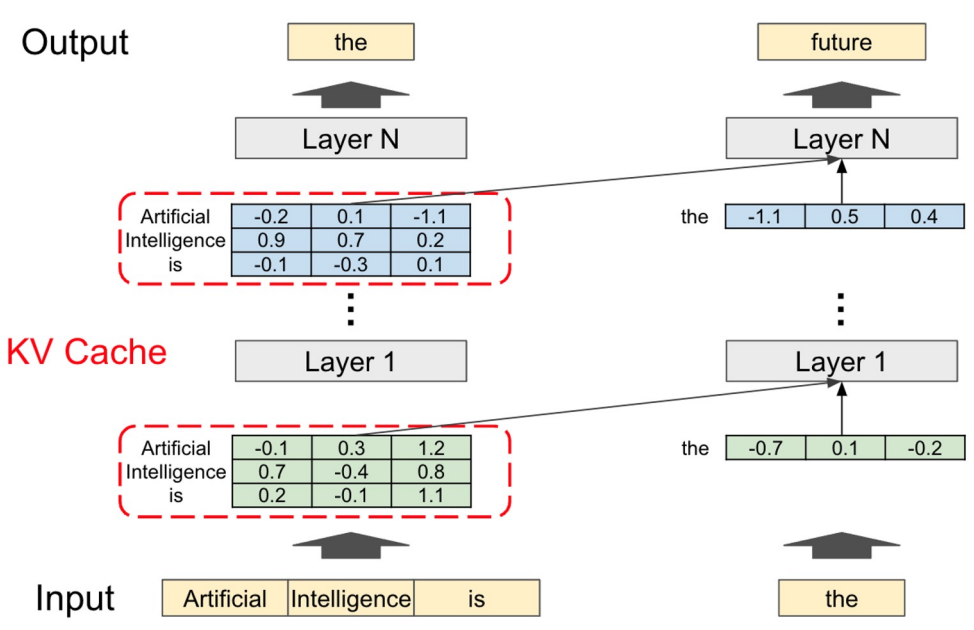
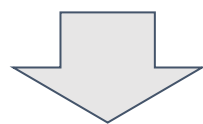
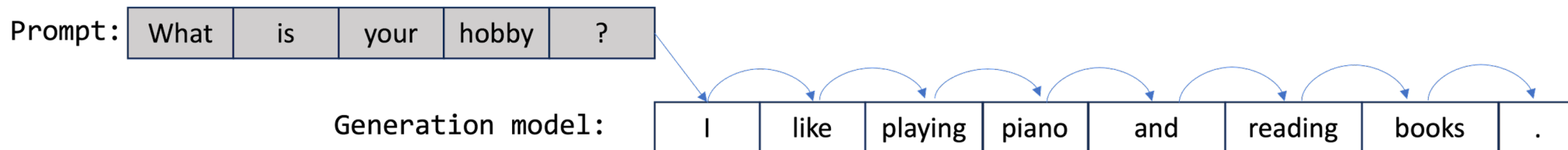


Table 1. Examples of neural network operations with their arithmetic intensities. Limiters assume FP16 data and an NVIDIA V100 GPU.

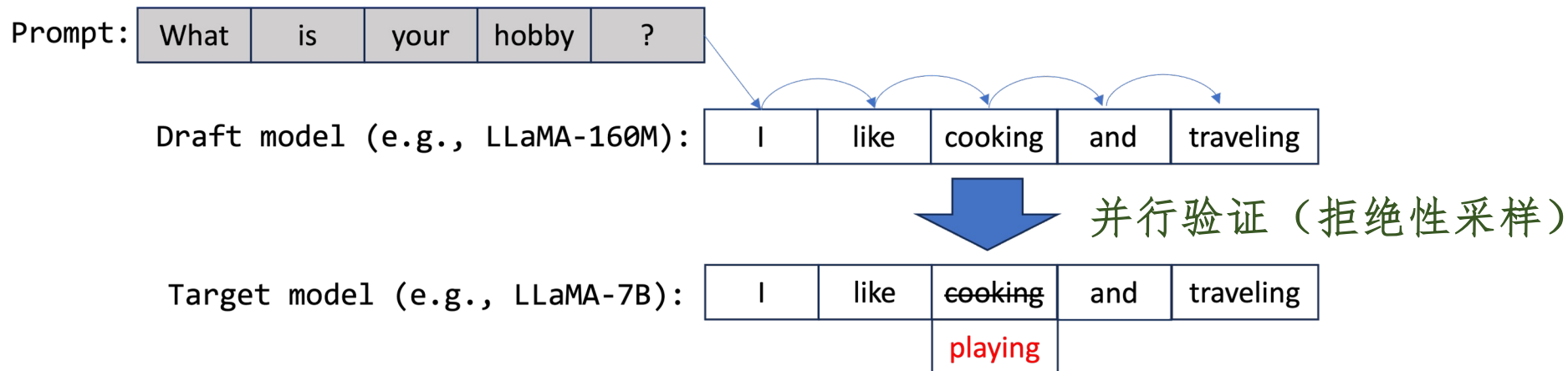
Operation	Arithmetic Intensity	Usually limited by...
Linear layer (4096 outputs, 1024 inputs, batch size 512)	315 FLOPS/B	arithmetic
Linear layer (4096 outputs, 1024 inputs, batch size 1)	1 FLOPS/B	memory
Max pooling with 3x3 window and unit stride	2.25 FLOPS/B	memory
ReLU activation	0.25 FLOPS/B	memory
Layer normalization	< 10 FLOPS/B	memory

# 并行推理方法：投机解码 (Speculative decoding, SD)

**KV cache**和模型权重的加载开销成倍下降

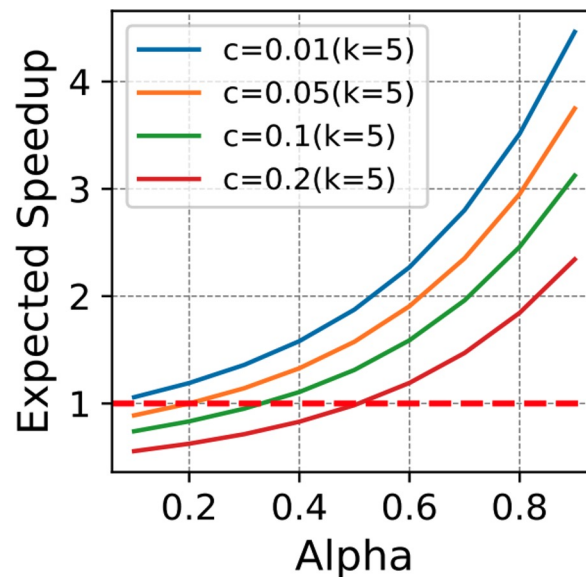
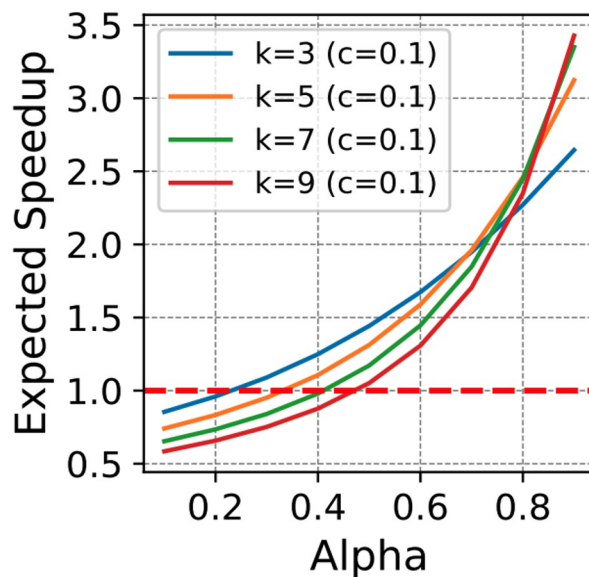


不是所有的token都“难”，可将大部分的顺序生成的负载转移到一个更小的草稿模型 (draft model)





# 投机解码的工作条件

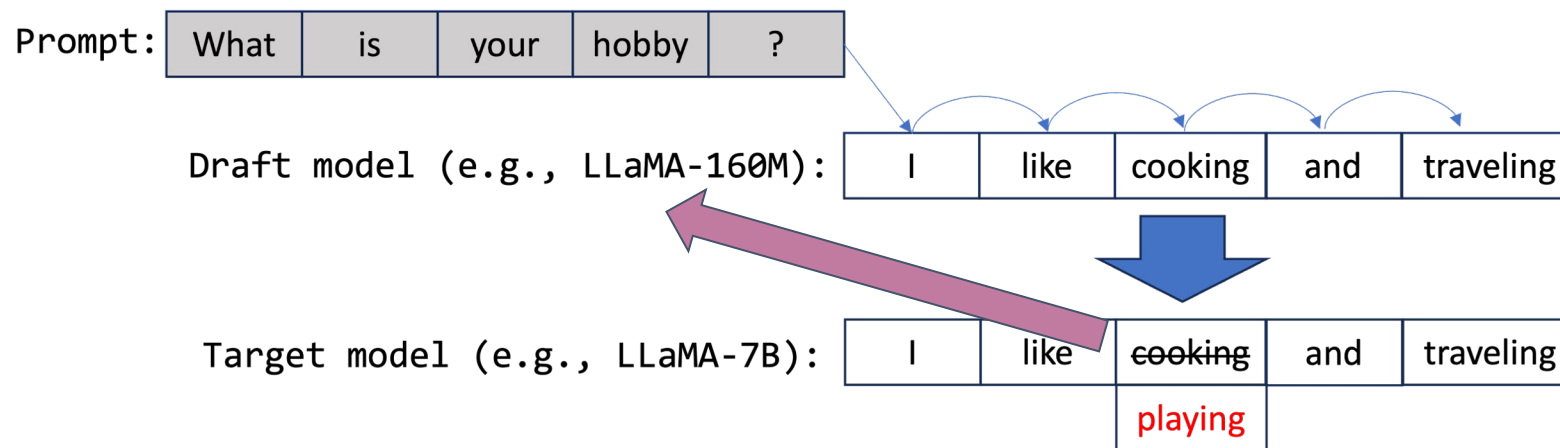


**c**: 草稿和目标模型的速度比；**k**: 每一步提议的token数；**Alpha**: token acceptance rate

- **结论：更高的token acceptance rate带来显著提高的加速效果**
  - 草稿模型必须足够逼近目标模型，同时保持较小的规模



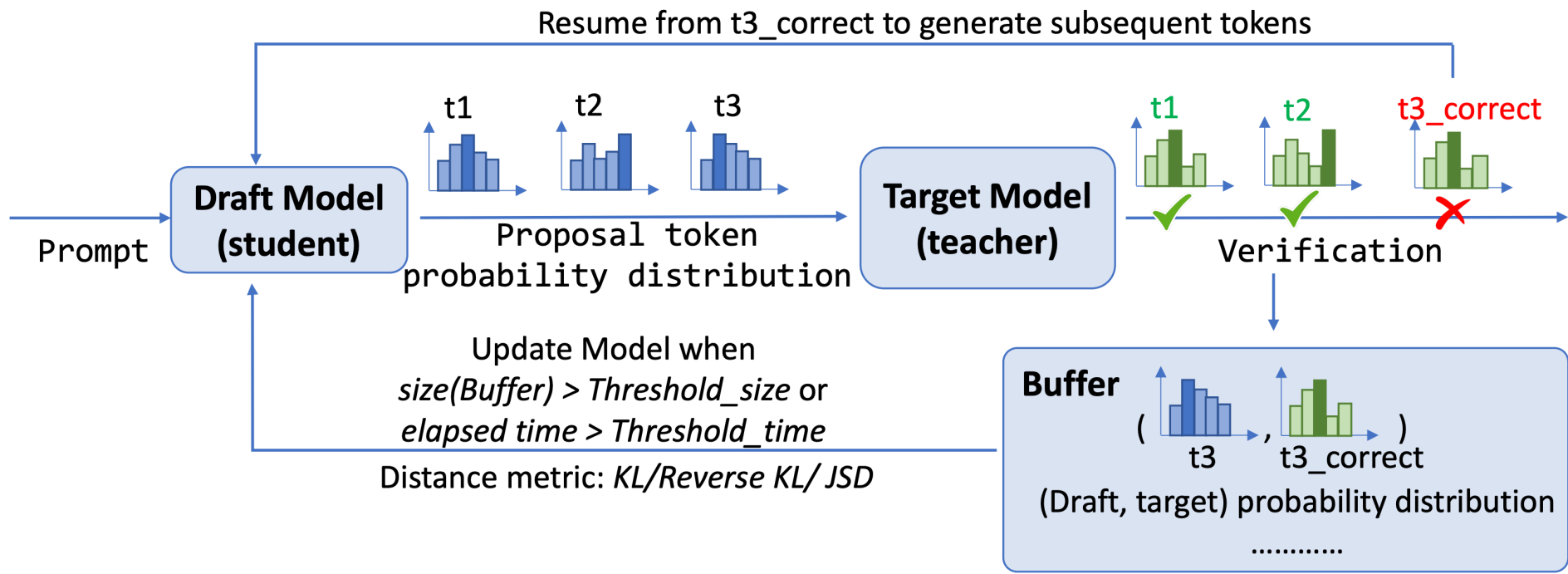
## 提高token acceptance rate的两个机会



- 投机解码可以检测草稿模型的错误并提供修正
  - 这些信息可以“免费”改进草稿模型，提高tokens acceptance rate，无需额外的标注成本
- 基于投机解码的LLM推理系统中的空闲算力（“spare” FLOPs）

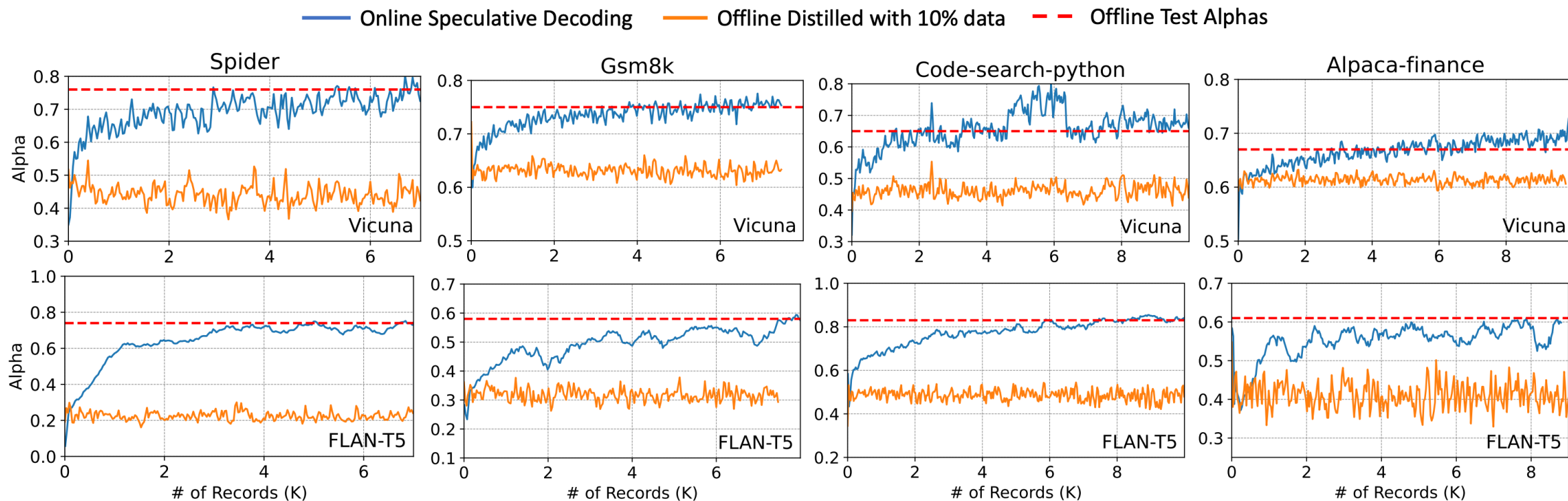
# 在线的投机解码 (OSD) : 在线蒸馏+投机解码

Open domain 情况下, draft model 快速适配 query dist.



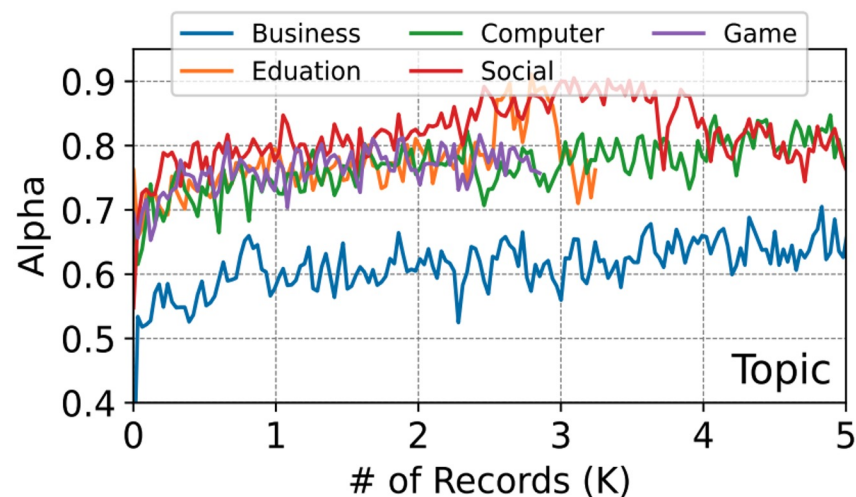
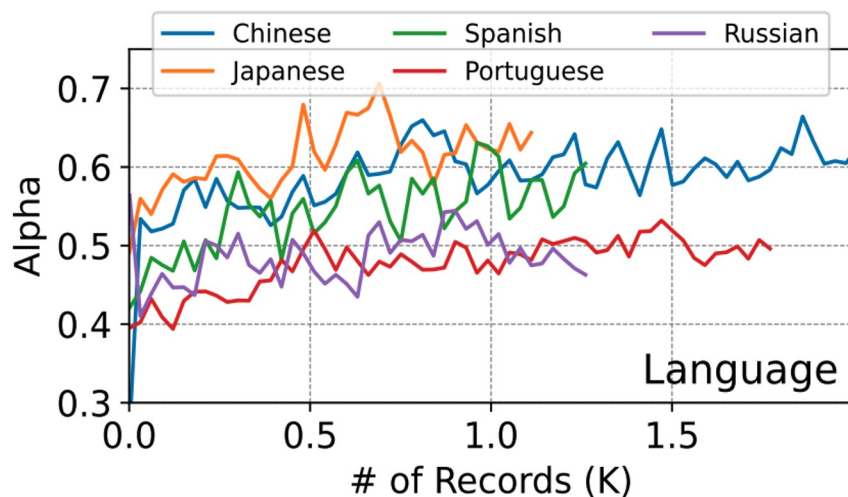
- 将草稿模型的错误预测和目标模型的校正结果存储在**buffer**中
- 当**buffer**打满, 基于在线蒸馏损失函数更新草稿模型

# 在线的投机解码：快速提高草稿模型的token acceptance rate



在线部署场景中，草稿模型渐渐适配数据分布，不断变准，因此加速效率不断提升

## 在线的投机解码：结合基于语言/主题的路由



- 使用多个草稿模型独立处理不同语言/主题，相对于单个草稿模型，进一步提高 **tokens acceptance rate**
- 可拓展为基于用户的路由，为每一个用户部署一个草稿模型



## 在线的投机解码：其他结果

Dataset	Spider	Chatbot Arena	Extra Parameters (B)
Medusa-7B	1.34×	2.03×	0.44
Medusa-7B + OSD	2.01×	2.38×	0.44
Draft model + OSD	2.17×	1.51×	0.16

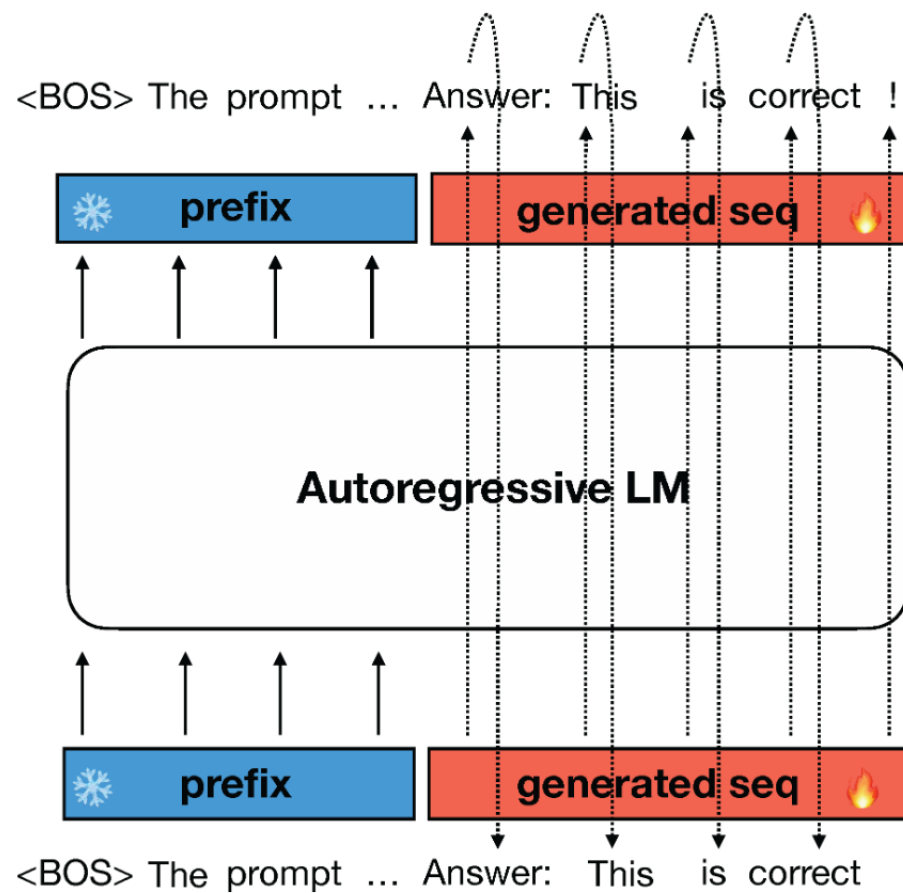
### 超越/结合Medusa

Dataset	Spider	Gsm8k	Alpaca-Finance	Code-Python
<b>Tokens with the greatest precision increase</b>	AV, SELECT, first, $\langle$ EOS $\rangle$ , template, SUM, G, COUNT, $\backslash$ n, city, WHERE, ', ;, (, IST, id	$\langle$ EOS $\rangle$ , >>, +, To, <<, this, =, %, know, are, We, calculate, be, The, have	1, Here, (, :, provide, depends, However, goals, amount, 3, there, The, $\backslash$ n, personal, will	'''', (, Here, python, ', how, doc, snippet, import, based, {, Python, This, :, you
<b>Tokens with the greatest recall increase</b>	SELECT, *, FROM, (, IST, *), $\backslash$ n, COUNT, G, first, WHERE, $\langle$ EOS $\rangle$ , IN, :, MAX, ', ;	start, >>, <<, +, find, how, we, =, fore, To, so, $\backslash$ , $\langle$ EOS $\rangle$ , then, let	general, 1, several, This, depends, Here, provide, However, goals, over, (, If, amount, it, can	Here, This, snippet, '''', ', how, python, (, takes, Python, you, doc, an, import, def

### Token acceptance rate 提升最多的tokens

# 语言模型自身仍是顺序解码器

- 自回归的学习方式
- 语言模型可以一次预测出多个tokens吗？



# Jacobi decoding: 将大语言模型变为并行解码器的初试

- 给定大语言模型  $p$ , 同时预测  $n$  个 token 等价于求解:

$$f(y_i, \mathbf{y}_{<i}, \mathbf{x}) = 0 \text{ for } i = 1, \dots, n.$$

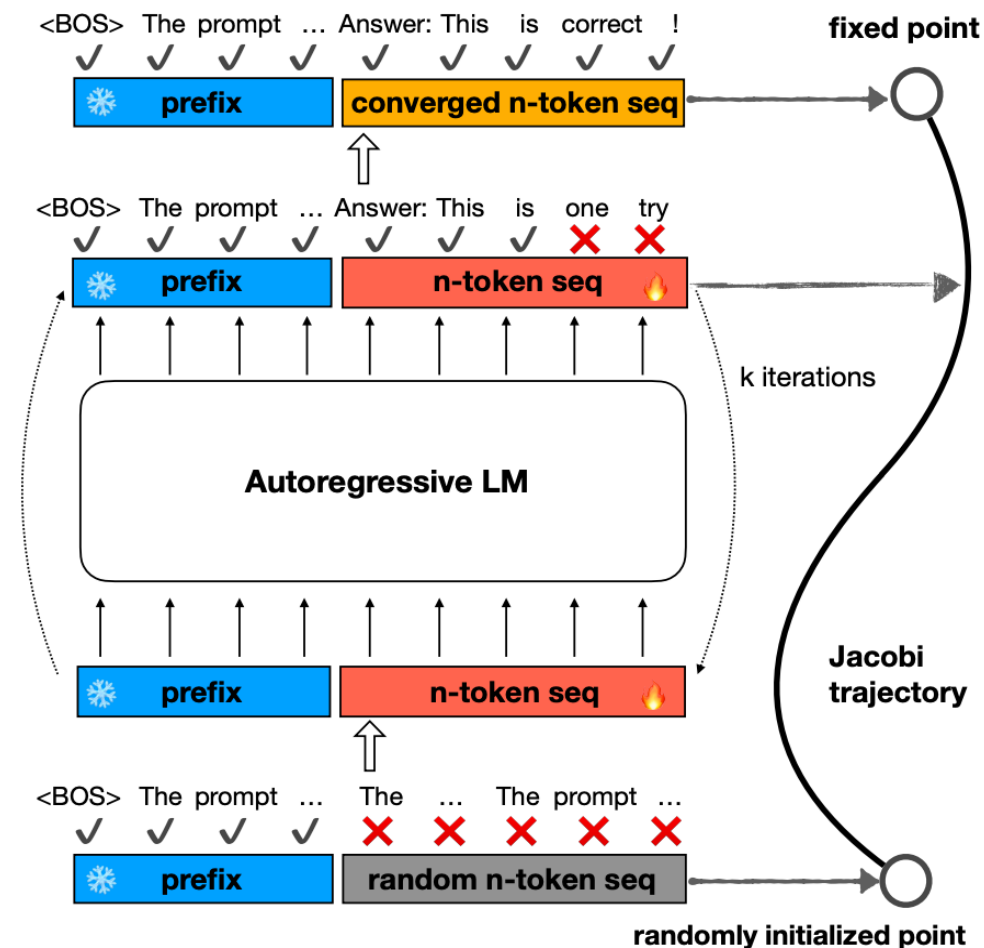
其中

$$f(y_i, \mathbf{y}_{<i}, \mathbf{x}) := y_i - \arg \max_y p(y | \mathbf{y}_{<i}, \mathbf{x})$$

- 这  $n$  个方程虽存在依赖, 但可被 **Jacobi** 不动点迭代法并行求解, 步数不超过  $n$ , 生成质量可保证

- 但实际效果差 (如: 仅 1.05 倍提升)

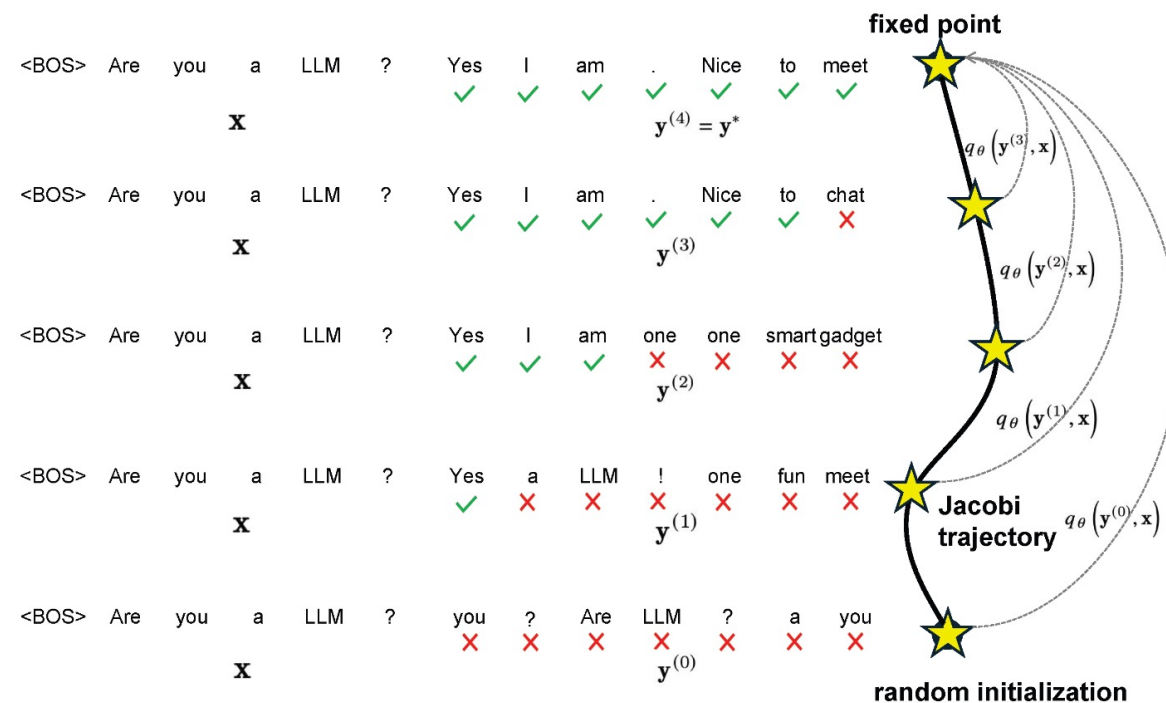
- 原因: 模型训练时未学过如何预测多个 tokens**





# 并行推理方法：一致性大语言模型 (Consistency LLMs, CLLMs)

- 如何习得预测n个tokens的能力？
  - 从随机初始化的起点预测**fixed point**？
    - 不行，问题太难，训练难收敛
  - 从**Jacobi**解码轨迹上的任意点预测**fixed point**？
    - 可以，形成一系列从简单到困难的学习问题，有助于模型收敛



# 一致性大语言模型：一致性损失函数

$$\mathcal{L}_{GC} = \mathbb{E}_{(\mathbf{x}, \mathcal{J}) \sim \mathcal{D}, \mathbf{y} \sim \mathcal{J}} \left[ \sum_{i=1}^n D(q_{\theta}(\cdot | \mathbf{y}_{:i}^*, \mathbf{x})) \| q_{\theta}(\cdot | \mathbf{y}_{:i}, \mathbf{x}) \right]$$

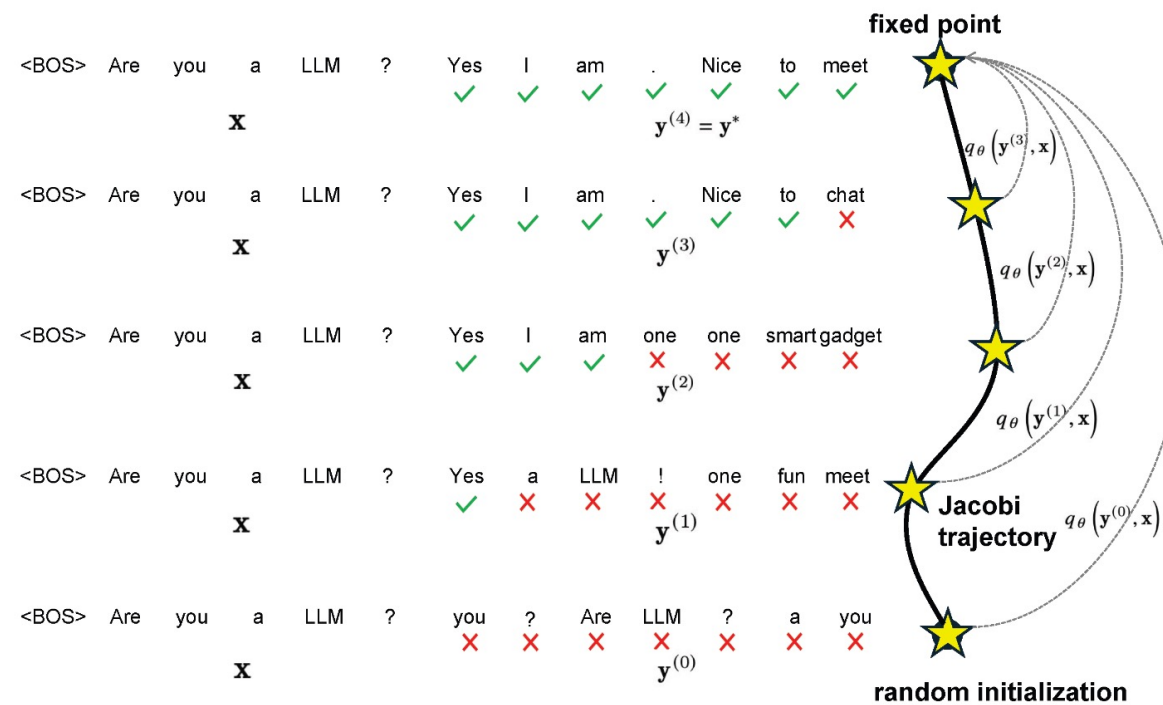
fixed point

$$\mathcal{L}_{LC} = \mathbb{E}_{(\mathbf{x}, \mathcal{J}) \sim \mathcal{D}, (\mathbf{y}^{(j)}, \mathbf{y}^{(j+1)}) \sim \mathcal{J}} \left[ \sum_{i=1}^n D(q_{\theta}(\cdot | \mathbf{y}_{:i}^{(j+1)}, \mathbf{x})) \| q_{\theta}(\cdot | \mathbf{y}_{:i}^{(j)}, \mathbf{x}) \right]$$

$D(\cdot \| \cdot)$  是分布间距离度量

$$\mathcal{L}_{AR} = \mathbb{E}_{(\mathbf{x}, \mathbf{l}) \sim \mathcal{D}} \left[ - \sum_{i=1}^N \log q_{\theta}(l_i | \mathbf{l}_{:i}, \mathbf{x}) \right]$$

自回归损失防止模型退化





# 一致性大语言模型：加速效果

```

USER: What are the most significant contributions Albert Einstein and Issac Newton each
have made to science and mathematics?
USER: What are the most significant contributions Albert Einstein and Issac Newton have
made to science and mathematics?

```

**VICUNA-7B**                      **CLLM-VICUNA-7B**

聊天， **Vicuna-7B**， **2.4倍加速**

```

USER: A tank has a capacity of 10000 gallons. Wanda and Ms. B decided to pump water from a pond to fi
ll the tank in two days. On the first day, working in shifts, Wanda filled 3/4 of the tank's capacity
with water, and Ms. B pumped 3/4 as much water as Wanda pumped into the tank that day. On the second
day, Wanda pumped 2/3 of the amount of water she pumped on the previous day, while Ms. B only pumped
1/3 of the number of gallons she pumped on the first day. How many gallons of water are remaining fo
r the tank to be full?
USER: A tank has a capacity of 10000 gallons. Wanda and Ms. B decided to pump water from a pond to fi
ll the tank in two days. On the first day, working in shifts, Wanda filled 3/4 of the tank's capacity
with water, and Ms. B pumped 3/4 as much water as Wanda pumped into the tank that day. On the second
day, Wanda pumped 2/3 of the amount of water she pumped on the previous day, while Ms. B only pumped
1/3 of the number of gallons she pumped on the first day. How many gallons of water are remaining fo
r the tank to be full?

```

**ABEL-7B-001**                      **CLLM-ABEL-7B-001**

数学， **Abel-7B-001**， **3倍加速**

```

USER: The SQL database has table named vehicle with columns ['Vehicle_ID', 'Model', 'B
uild_Year', 'Top_Speed', 'Power', 'Builder', 'Total_Production'], table named driver w
ith columns ['Driver_ID', 'Name', 'Citizenship', 'Racing_Series'], table named vehicle
_driver with columns ['Driver_ID', 'Vehicle_ID'], Question: What are the vehicle ids a
nd models which have been driven by more than 2 drivers or been driven by the driver n
amed 'Jeff Gordon'?
USER: The SQL database has table named vehicle with columns ['Vehicle_ID', 'Model', 'B
uild_Year', 'Top_Speed', 'Power', 'Builder', 'Total_Production'], table named driver w
ith columns ['Driver_ID', 'Name', 'Citizenship', 'Racing_Series'], table named vehicle
_driver with columns ['Driver_ID', 'Vehicle_ID'], Question: What are the vehicle ids a
nd models which have been driven by more than 2 drivers or been driven by the driver n
amed 'Jeff Gordon'?

```

**Deepseek-Coder-7B**                      **CLLM-Deepseek-Coder-7B**

代码， **Deepseek-coder-7B**， **3.4倍加速**



# 一致性大语言模型：系统性比较

- 至多**3.6**倍加速
- 相比于**Medusa2**，不需要模型架构上的改变
- 有保证的生成质量

Methods	Speed (tokens/s)	Speedup	Metric	Size
GSM8K				
Fine-tuned LLaMA2-7B (Chern et al.)				
+ AR	43.5	1.0×	59.1	6.7B
+ Jacobi	45.7	1.1×	59.1	
+ lookahead	74.8	1.7×	59.1	
CLLM-LLaMA2-7B				
+ AR	43.5	1.0×	56.4	6.7B
+ Jacobi	132.4	<b>3.0×</b>	56.4	
+ lookahead	125.2	<b>2.9×</b>	56.4	
Medusa-2 + LLaMA2-7B				
+ typical	70.2	1.6×	51.3	8.3B
Fine-tuned LLaMA2-7B + distilled LLaMA-160m				
+ speculative	73.8	1.7×	59.1	6.8B
ShareGPT (MT-Bench)				
Fine-tuned LLaMA2-7B				
+ AR	37.6	1.0×	6.5	6.7B
+ Jacobi	39.9	1.1×	6.5	
+ lookahead	60.8	1.6×	6.5	
CLLM-LLaMA2-7B				
+ AR	36.7	1.0×	6.4	6.7B
+ Jacobi	88.4	<b>2.4×</b>	6.4	
+ lookahead	95.0	<b>2.5×</b>	6.4	
Medusa-2 + LLaMA2-7B				
+ typical	102.5	2.7×	6.4	8.3B
Fine-tuned LLaMA2-7B + distilled LLaMA-160m				
+ speculative	51.3	1.4×	6.5	6.8B

Methods	Speed (tokens/s)	Speedup	Metric	Size
Spider				
Fine-tuned Deepseek-7B				
+ AR	38.0	1.0×	70.0	6.7B
+ Jacobi	39.5	1.0×	70.0	
+ lookahead	55.3	1.5×	70.0	
CLLM-Deepseek-7B				
+ AR	38.0	1.0×	69.3	6.7B
+ Jacobi	127.4	<b>3.4×</b>	69.3	
+ lookahead	135.2	<b>3.6×</b>	69.3	
Medusa-2 + Deepseek-7B				
+ typical	104.2	2.7×	66.4	8.3B
Fine-tuned Deepseek-7B + distilled LLaMA-160m				
+ speculative	66.8	1.8×	70.0	6.8B
Code-Search-Net Python				
Fine-tuned Deepseek-7B				
+ AR	40.1	1.0×	60.4	6.7B
+ Jacobi	43.2	1.1×	60.4	
+ lookahead	68.0	1.7×	60.0	
CLLM-Deepseek-7B				
+ AR	38.5	1.0×	59.2	6.7B
+ Jacobi	102.1	<b>2.5×</b>	59.2	
+ lookahead	115.7	<b>2.9×</b>	59.2	
Medusa-2 + Deepseek-7B				
+ typical	128.0	3.2×	48.3	8.3B
Fine-tuned Deepseek-7B + distilled LLaMA-160m				
+ speculative	59.3	1.5×	60.4	6.8B



# 一致性大语言模型：训练的开销很低

Table 9. Computation required for consistency training.

Dataset	Training time	% of pre-training cost	Training resources
Spider	2 hours	$< 0.01\%$	8 A100 40GB GPUs
GSM8K	12 hours	$\sim 0.01\%$	8 A100 40GB GPUs
CodeSearchNet-Python	22 hours	$\sim 0.1\%$	8 A100 40GB GPUs
ShareGPT	30 hours	$\sim 0.2\%$	8 A100 40GB GPUs



# 一致性大语言模型：加速根源

## Target LLM

```

0: with ## __ ' , Manager table have number ID ' Could Player or ruction ['
1: SELECT T as ## as '' , as columns _ _ C l you _ ''
2: SELECT country 2 T FROM ## Country COUNT ## FROM name name ub _ t FROM
3: SELECT country FROM FROM 1 player T AS ( WHERE player AS GROUP WHERE id WHERE
4: SELECT country FROM player player WHERE GROUP 1 T SELECT T . T BY player AS
5: SELECT country FROM player GROUP _ country _ . 1 country GROUP country GROUP COUNT GROUP
6: SELECT country FROM player GROUP BY CON H count H H H BY H BY (
7: SELECT country FROM player GROUP BY country TRY I L H AV AV country H H
8: SELECT country FROM player GROUP BY country H A I A V AV ING G AV
9: SELECT country FROM player GROUP BY country H AV ING V ING COUNT COUNT COUNT >
10: SELECT country FROM player GROUP BY country H AV ING count E _ ( ( (
11: SELECT country FROM player GROUP BY country H AV ING count ( ( number *) *)
12: SELECT country FROM player GROUP BY country H AV ING count ( *) *) *) >
13: SELECT country FROM player GROUP BY country H AV ING count ( *) > *) >
14: SELECT country FROM player GROUP BY country H AV ING count ( *) > 1 'ln'
  
```

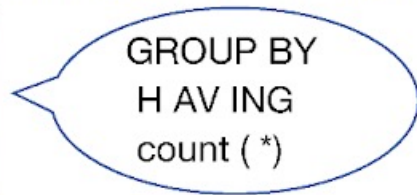
## Consistency LLM

```

0: with ## __ ' , Manager table have number ID ' Could Player or ruction ['
1: SELECT country country FROM player GROUP BY country H H AV AV AV ING *) *)
2: SELECT country FROM player BY BY H AV AV AV ING count ( *) > ''
3: SELECT country FROM player GROUP BY H AV ING count ( *) > 1 '' ''
4: SELECT country FROM player GROUP BY country AV AV count ( *) > 1 '' 'ln'
5: SELECT country FROM player GROUP BY country H AV ING count *) > 1 '' 'ln'
6: SELECT country FROM player GROUP BY country H AV ING count ( *) > 1 'ln'
  
```



Target LLM



a lot of collocations



CLLM

- **Fast forwarding:** 多个连续的tokens在单次网络传播中被正确预测
- **Stationary tokens:** 提前被正确预测，并在后续迭代中保持不变，即使之前的tokens存在错误

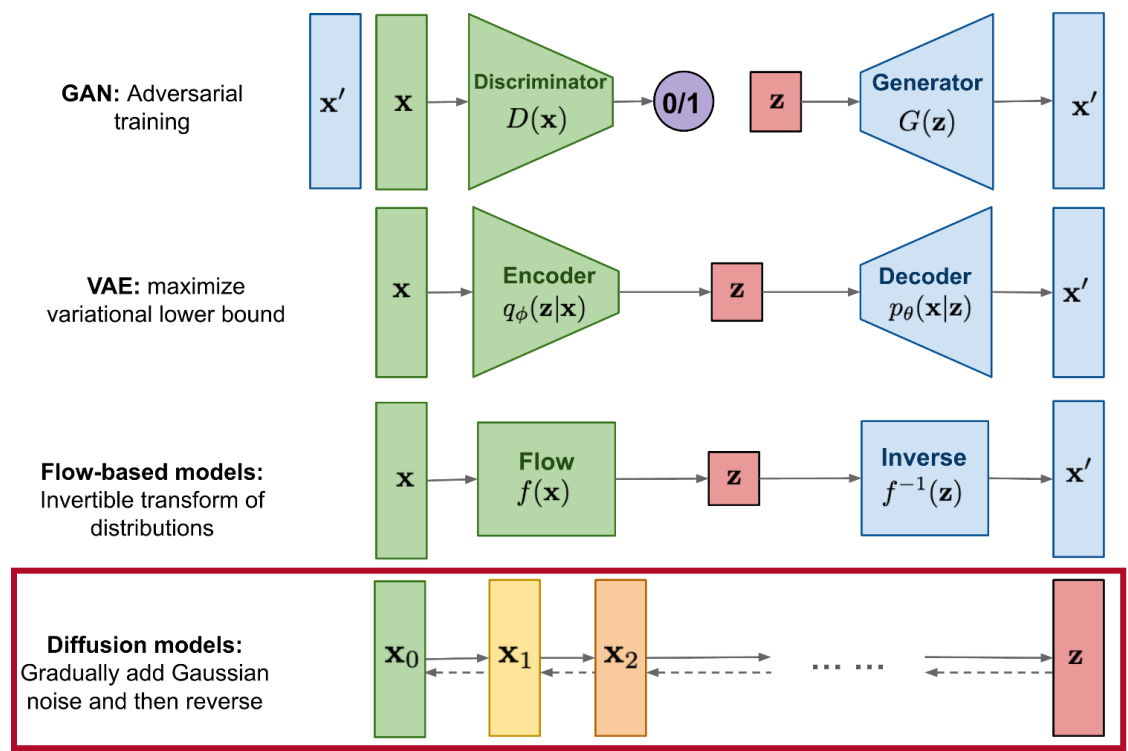
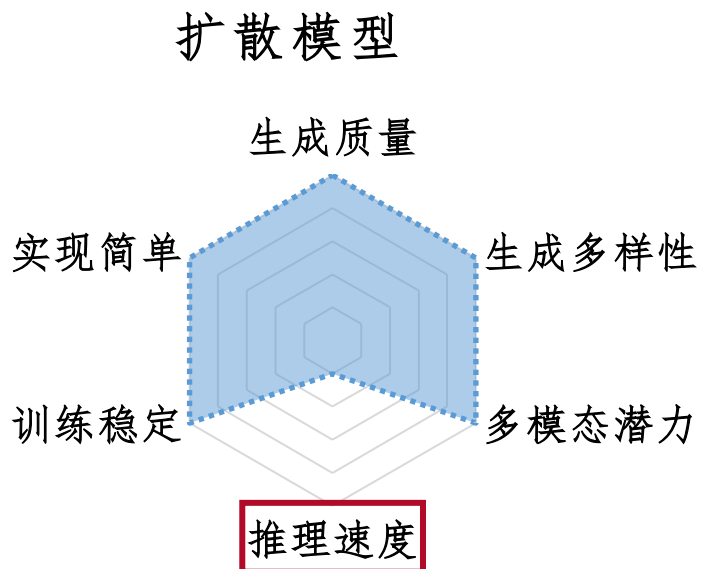


# AIGC大模型的高效推理方法

- 大语言模型的并行推理方法（15min）
  - 投机解码
  - 一致性蒸馏
- 大扩散模型的低步数推理方法（15min）
  - 采样器设计
  - 一致性蒸馏
- 大模型架构、序列状态、缓存等方面的优化方法（5min）

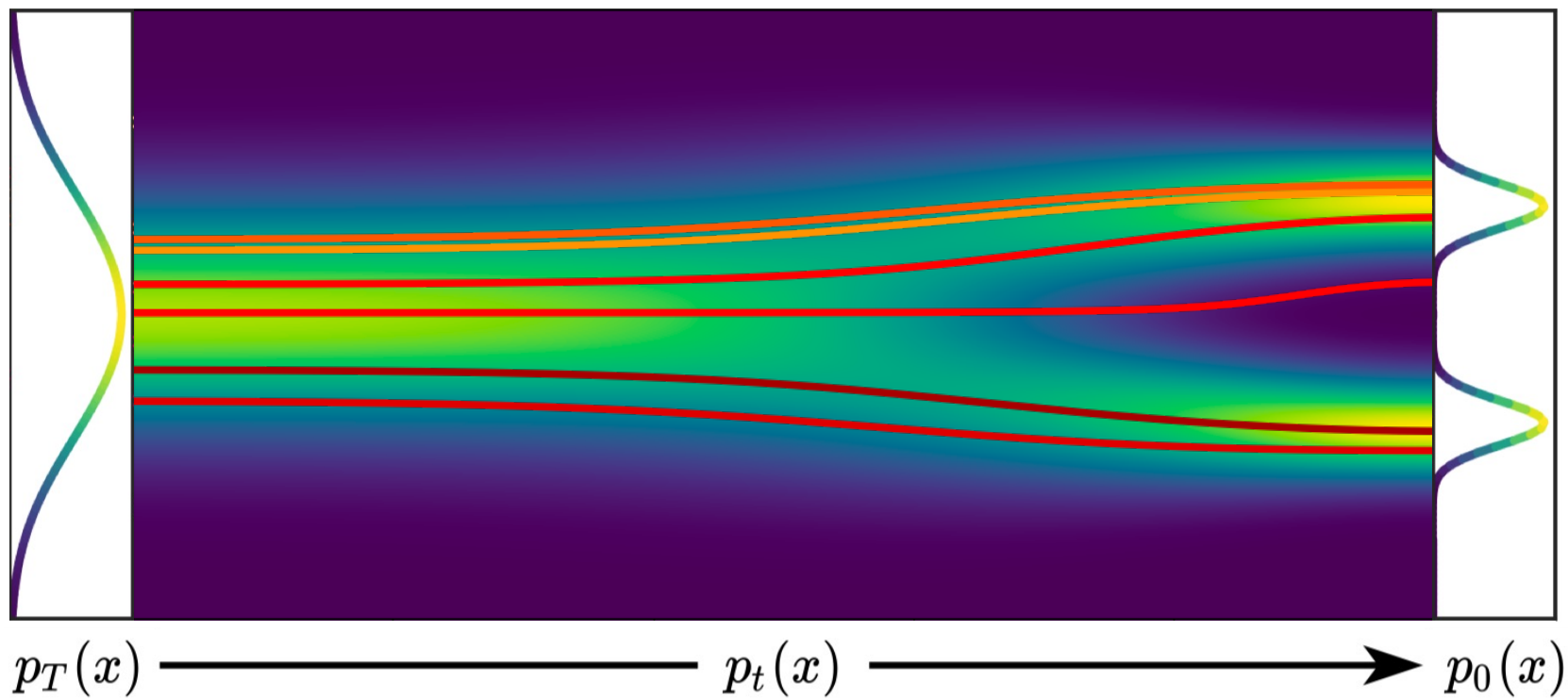


# 采样效率是扩散模型的瓶颈问题：每个去噪步都需要一次模型前传



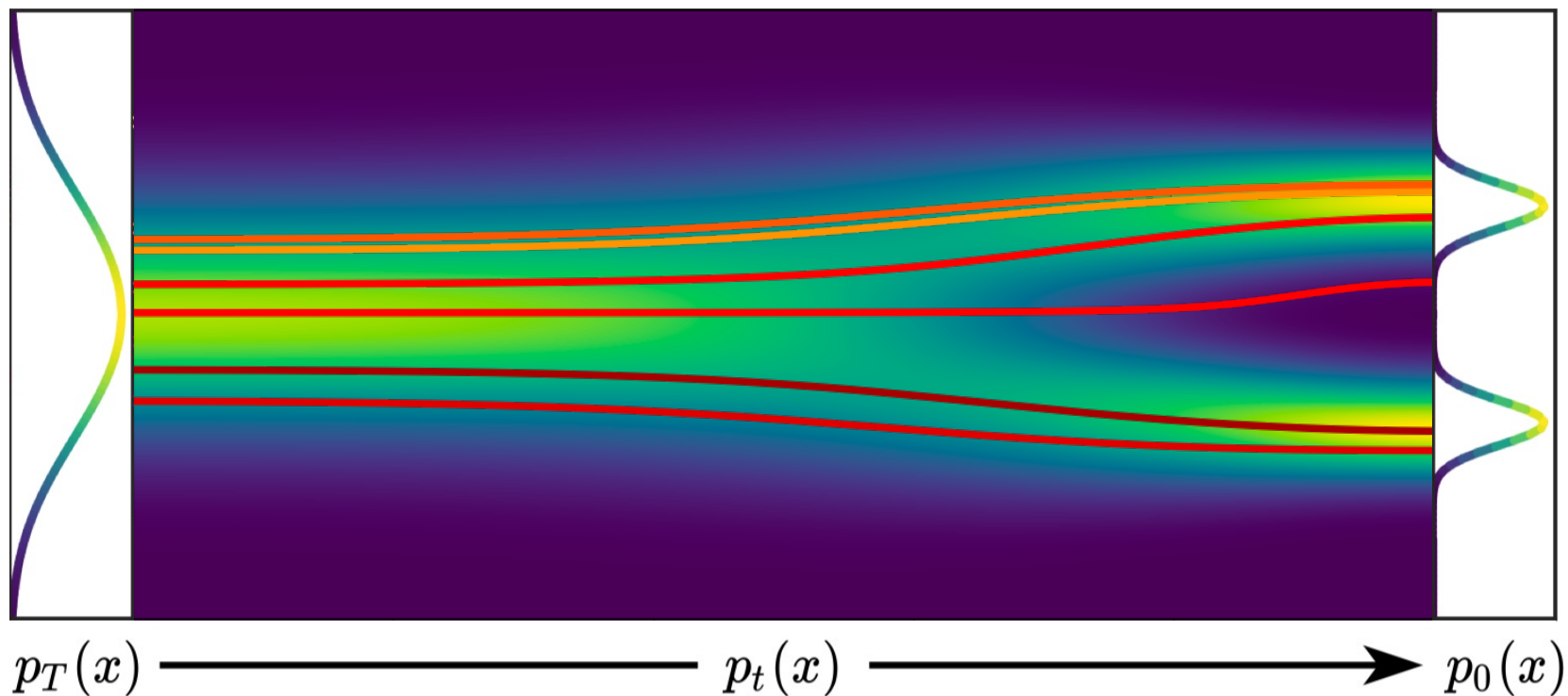
原则上需要迭代很多次，每次需调用神经网络

# 出发点：常微分方程的视角



时时刻刻边缘分布相同，但是路径上没有噪声！

# 扩散模型采样等价常微分方程离散化



Sampling method	Steps to converge
Traditional <b>SDE</b> Solvers	$\sim 200$
Traditional <b>ODE</b> Solvers	$\sim 100$

Checkpoint: CIFAR-10, VP

[Song et al., ICLR'21]

# DPM-Solver: 面向扩散概率模型的常微分方程离散化

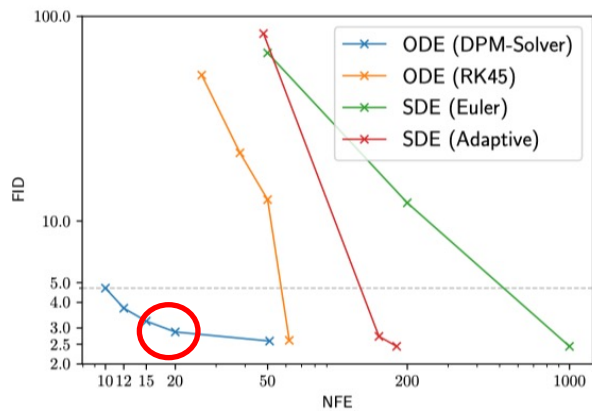
- 针对扩散概率模型半线性等特点，设计等价常微分方程的离散化解析形式

经典龙格库塔法  $\mathbf{x}_t = \mathbf{x}_s + \int_s^t \left( f(\tau)\mathbf{x}_\tau + \frac{g^2(\tau)}{2\sigma_\tau} \epsilon_\theta(\mathbf{x}_\tau, \tau) \right) d\tau$  整体黑盒泰勒展开并做差分近似

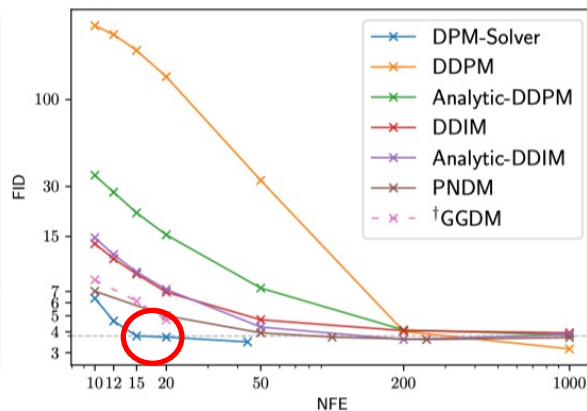
所提 DPM-Solver  $\mathbf{x}_{t_{i-1} \rightarrow t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \sum_{n=0}^{k-1} \hat{\epsilon}_\theta^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}}) \int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda + \mathcal{O}(h_i^{k+1})$

解析形式
神经网络部分差分近似
解析形式
高阶小量

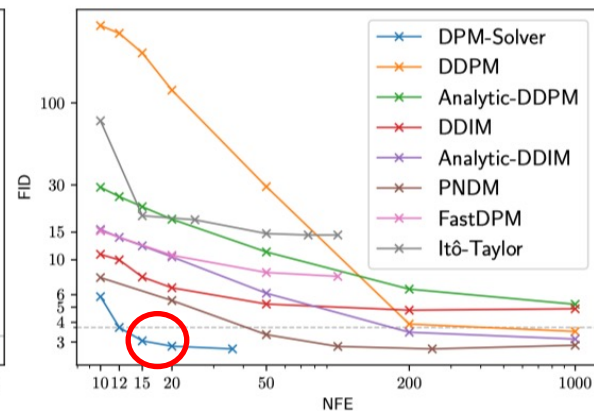
# DPM-Solver: 结果



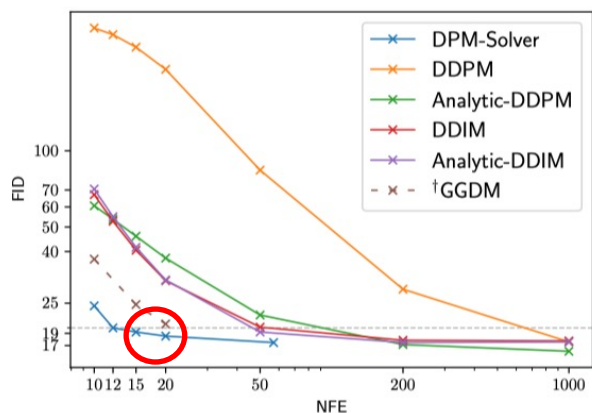
(a) CIFAR-10 (continuous)



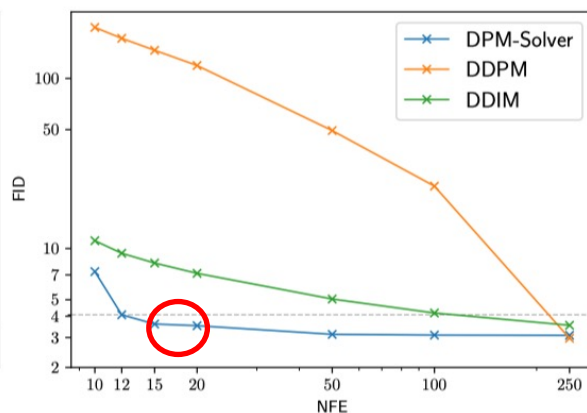
(b) CIFAR-10 (discrete)



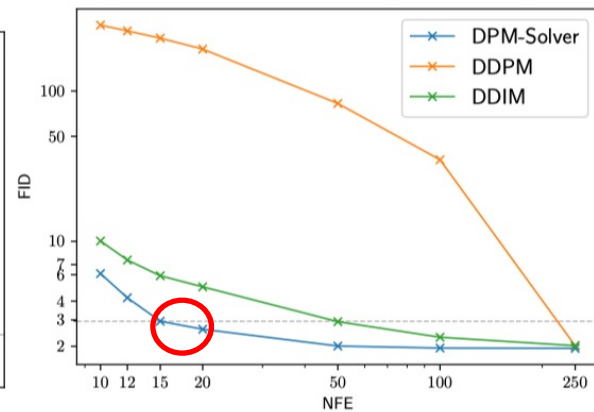
(c) CelebA 64x64 (discrete)



(d) ImageNet 64x64 (discrete)



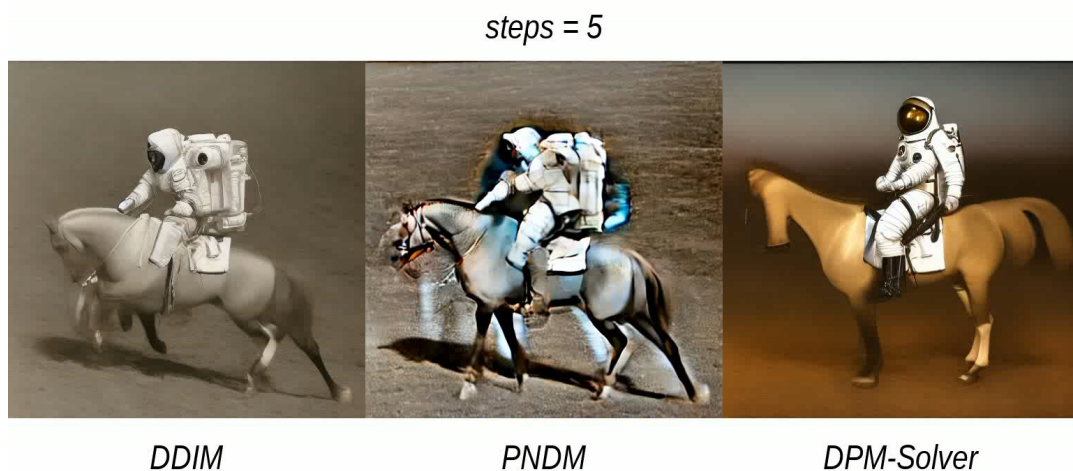
(e) ImageNet 128x128 (discrete)



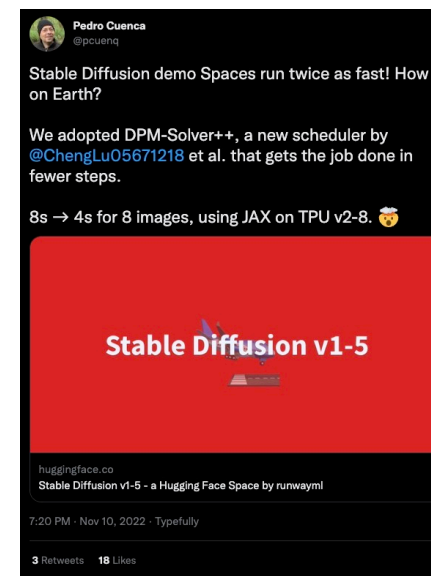
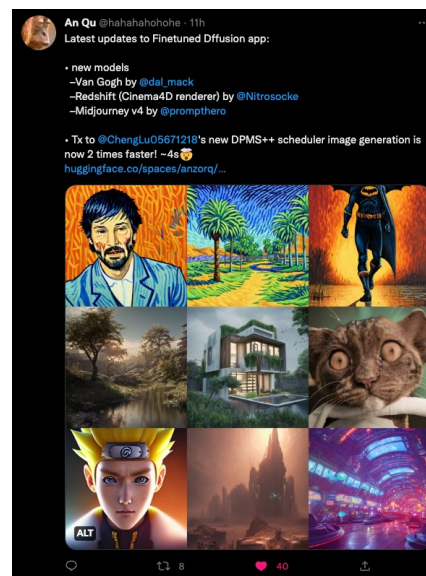
(f) LSUN bedroom 256x256 (discrete)

# DPM-Solver: 结果

- 是当前最快的无需额外学习的扩散概率模型采样算法，**15步生成高清图像**
- 被多个主流开源社区（**Github** 累计星标**6万余次**）支持/设为默认算法

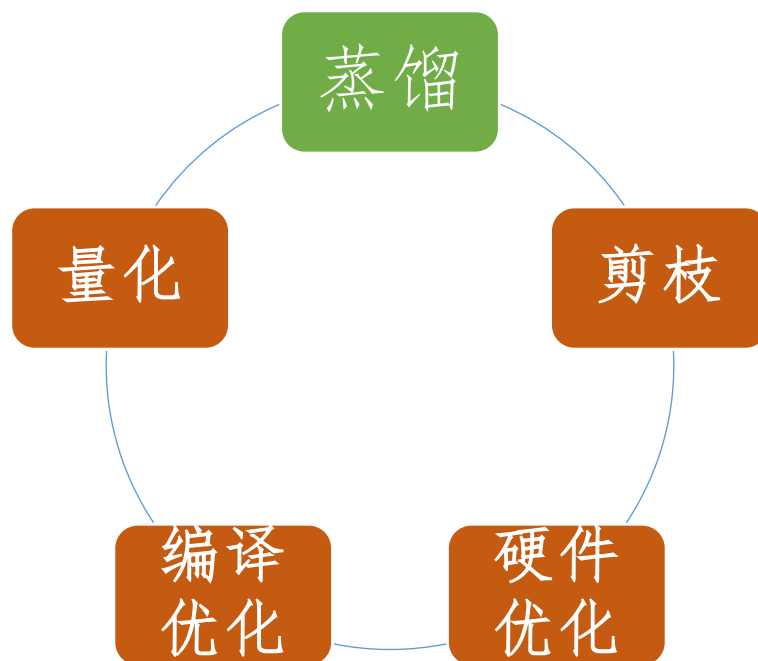


根据文本输入 15 步生成 512x512 高清图像



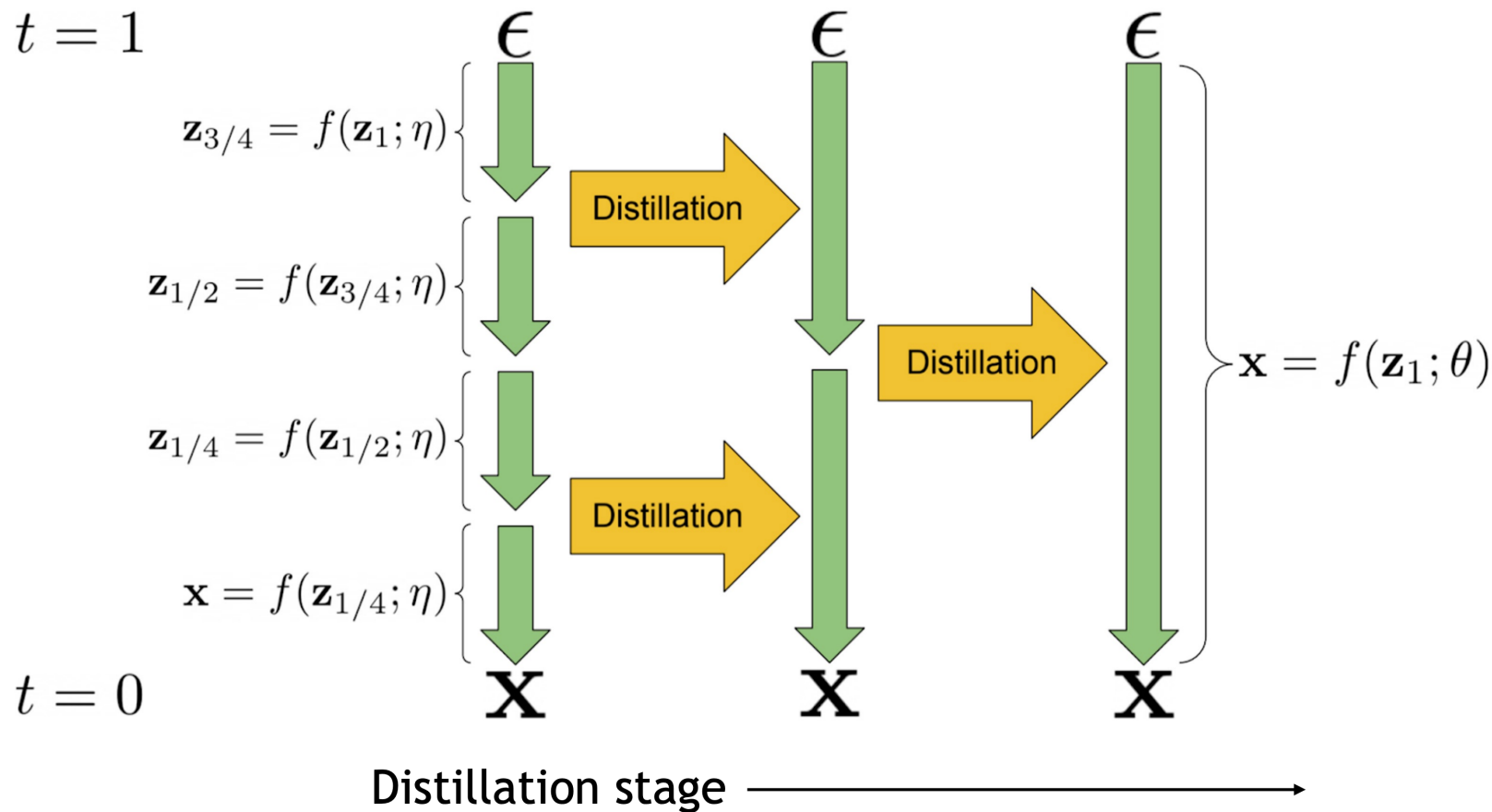
著名开源模型 **Stable Diffusion** 等官方宣传

## 从模型蒸馏的角度出发实现低步数推理





# 渐进蒸馏



- 提前确定noise schedule
- 基于预训练模型，以**确定性 ODE sampler**作为教师
- **2步蒸馏为1步**：学生只需要教师的一半采样步
- 循环往复

# 改进： On Distillation of Guided Diffusion Models (Meng et al., CVPR 2023)

- 带来了：
  - CF-Guidance
  - Stochastic sampling
  - Text-to-image
  - Image-to-image
  - Inpainting
  - Latent Diffusion



问题：蒸馏多个模型，训练成本高

1-4步即可生成高质量大图

# 一致性蒸馏

- 概率流ODE定义了从噪声到数据的一一映射

=>直接建模此映射

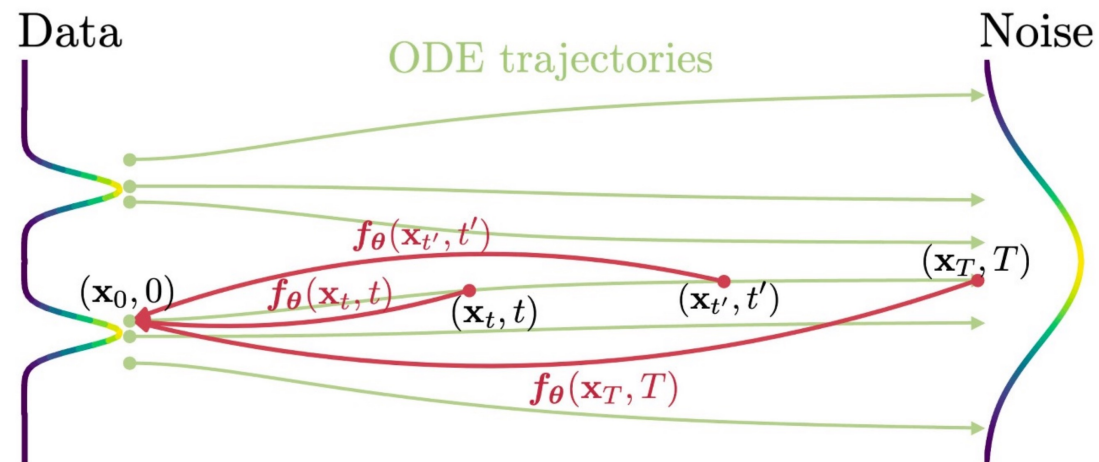
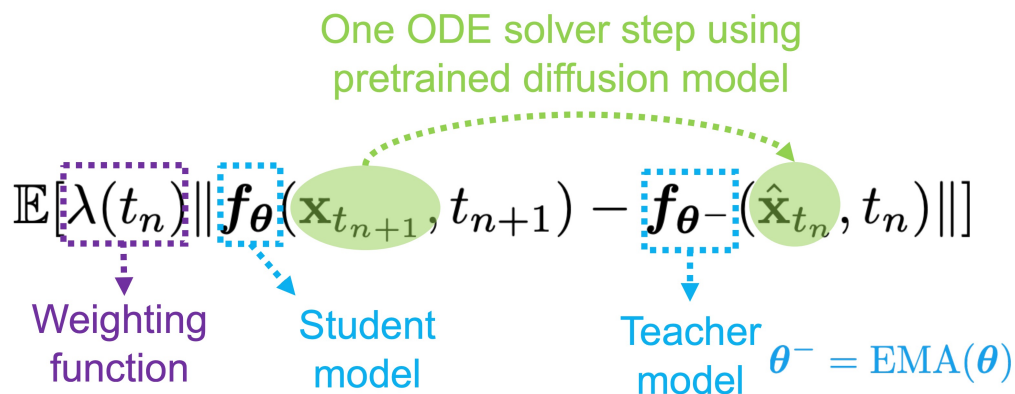
$$\forall t \in [0, T] : f_{\theta}(\mathbf{x}_t, t) = \mathbf{x}_0$$

- 模型参数化: 需保证边界条件 ( $t=0$ )

$$f_{\theta}(\mathbf{x}, t) = c_{\text{skip}}(t)\mathbf{x} + c_{\text{out}}(t)F_{\theta}(\mathbf{x}, t)$$

$$c_{\text{skip}}(0) = 1 \quad c_{\text{out}}(0) = 0$$

- 训练:



## Algorithm 1 Multistep Consistency Sampling

**Input:** Consistency model  $f_{\theta}(\cdot, \cdot)$ , sequence of time points  $\tau_1 > \tau_2 > \dots > \tau_{N-1}$ , initial noise  $\hat{\mathbf{x}}_T$

$\mathbf{x} \leftarrow f_{\theta}(\hat{\mathbf{x}}_T, T)$

**for**  $n = 1$  **to**  $N - 1$  **do**

    Sample  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$\hat{\mathbf{x}}_{\tau_n} \leftarrow \mathbf{x} + \sqrt{\tau_n^2 - \epsilon^2} \mathbf{z}$

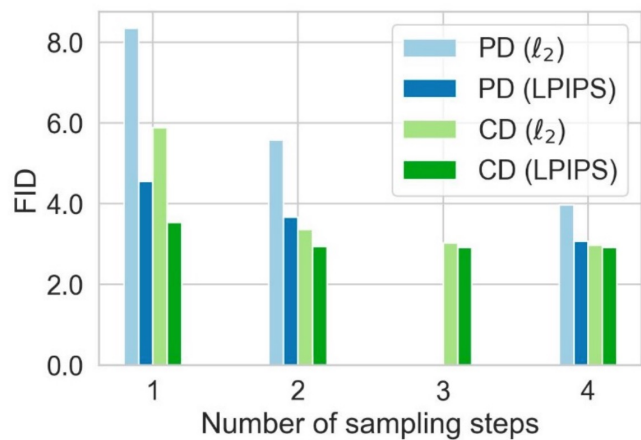
$\mathbf{x} \leftarrow f_{\theta}(\hat{\mathbf{x}}_{\tau_n}, \tau_n)$

**end for**

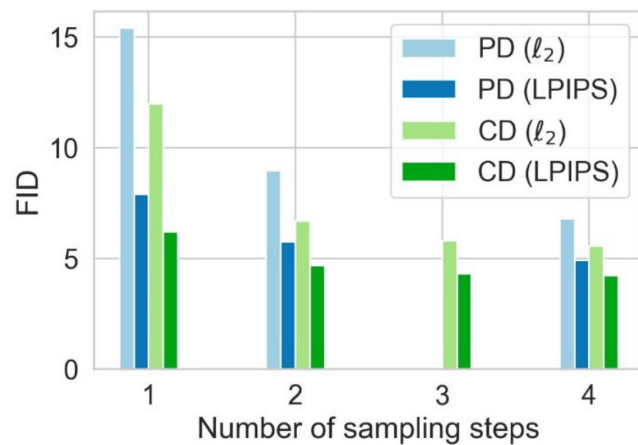
**Output:**  $\mathbf{x}$



# 一致性蒸馏：结果



(a) CIFAR-10

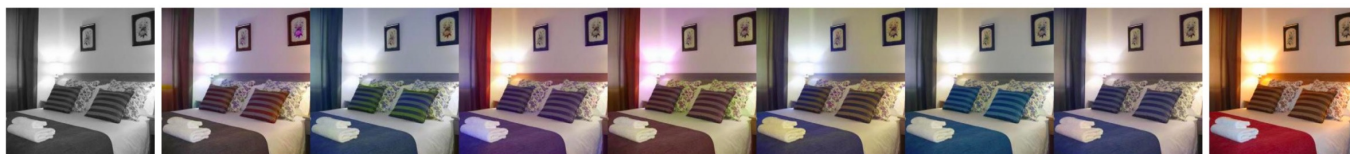


(b) ImageNet 64 × 64

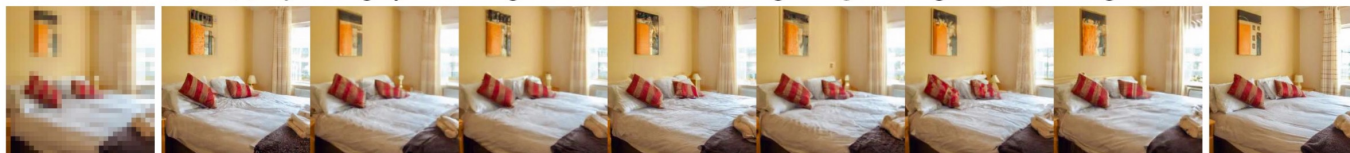
一步生成的 **SOTA** FID:

- 3.55 on CIFAR-10
- 6.20 on ImageNet 64

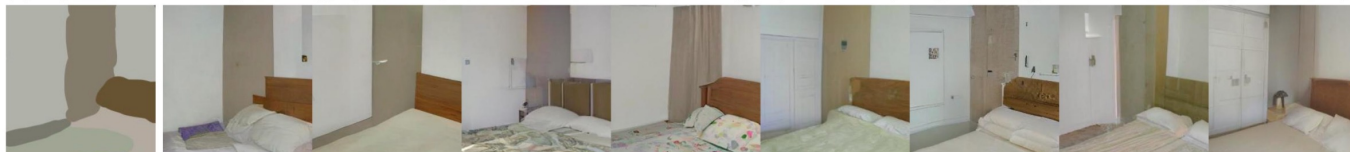
可用于 **零样本图  
像编辑** 等应用



(a) *Left*: The gray-scale image. *Middle*: Colorized images. *Right*: The ground-truth image.



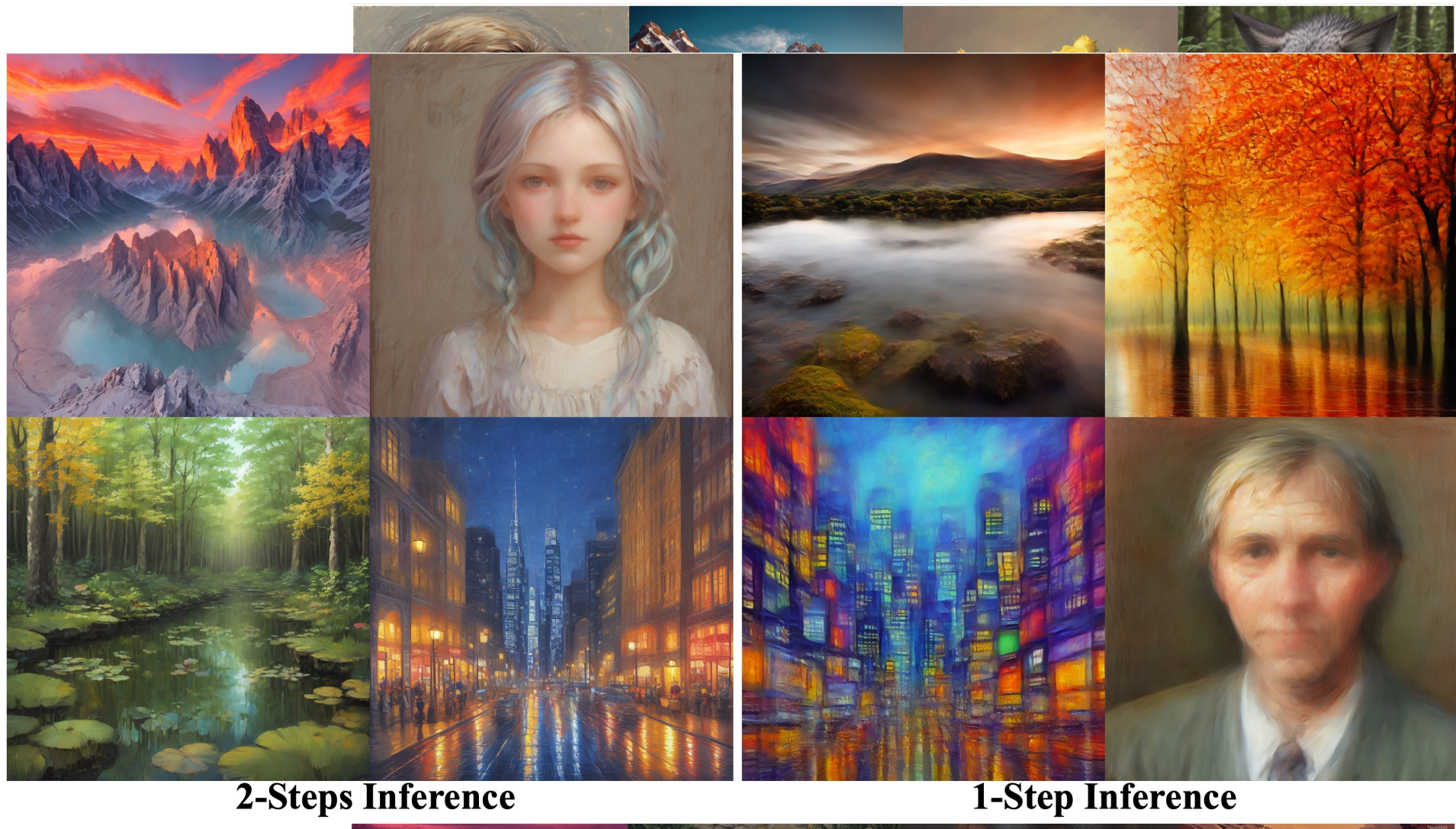
(b) *Left*: The downsampled image (32 × 32). *Middle*: Full resolution images (256 × 256). *Right*: The ground-truth image (256 × 256).



(c) *Left*: A stroke input provided by users. *Right*: Stroke-guided image generation.



# 隐空间一致性蒸馏



2-Steps Inference

1-Step Inference

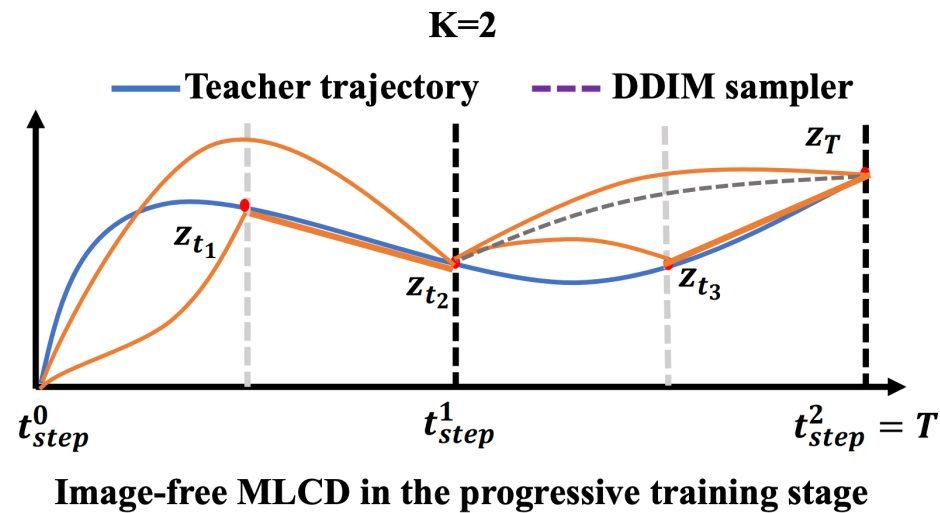
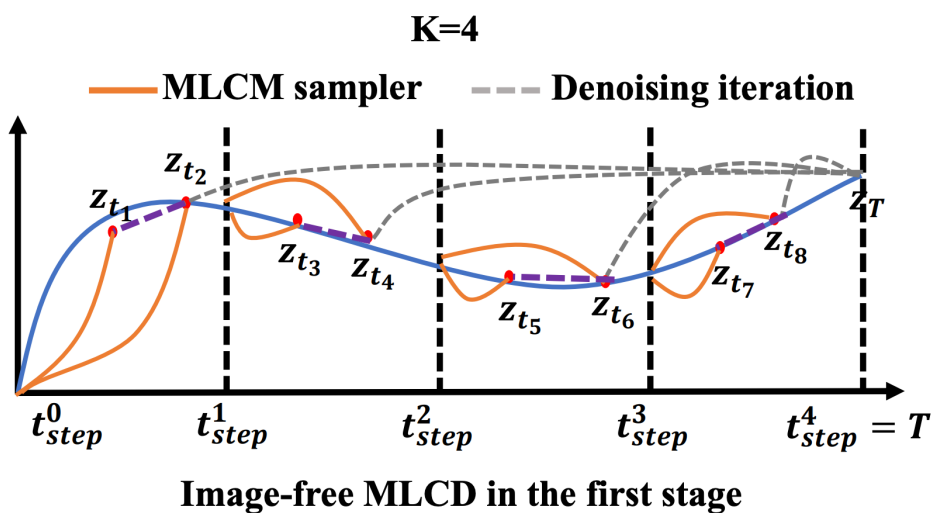
4-Steps Inference

从预训练SD蒸馏:

- 4,000个训练步  
(约32个A100 GPU hour)
- 生成高质量的  
768×768分辨率  
图像

# 多步隐空间一致性蒸馏

- 直接学习从噪声到图像的一致性映射太难，可将**ODE**轨迹切段分别处理
  - 课程学习思想：渐进减少切段数目
  - 基于**教师模型**的采样进行学习，避免对高质量训练集的依赖
  - 引入**偏好学习**损失，兼顾加速和对齐





# 多步隐空间一致性蒸馏：结果



Figure 1:  $1024 \times 1024$  image samples from MLCM, distilled from SDXL-base-1.0 [32] based on LoRA [8]. From top to bottom, 2, 3, and 4 sampling steps are adopted, respectively. Apart from such good visual quality, MLCM can also yield improved metrics compared to strong baselines.

Table 1: Quantitative comparisons on MSCOCO-2017 5K validation datasets. All models adopts SDXL architecture.

Method	Step	CS	AS	IR	PS
DDIM [42]	25	33.36	5.54	0.87	0.229
LCM [23]	4	32.53	5.42	0.48	0.224
SDXL-Turbo [38]	4	33.30	5.64	0.83	0.226
SDXL-Lightning [17]	4	32.40	5.63	0.72	0.229
SDXL-Lightning [17]	8	32.73	5.95	0.71	0.227
HyperSD [33]	4	32.64	5.52	1.15	<b>0.234</b>
HyperSD [33]	8	32.41	5.83	1.14	0.233
MLCM	2	33.02	6.10	1.10	0.227
MLCM	3	33.24	6.17	1.18	0.232
MLCM	4	33.30	6.19	1.20	0.233
MLCM	8	<b>33.48</b>	<b>6.20</b>	<b>1.22</b>	0.233

- 无需训练图片，统一模型，兼容不同采样步数
- 2步采样即可实现高质量的 $1024^2$ 的图像生成，远超SDXL, HyperSD等开源模型
- 赋能OPPO产品



# 多步隐空间一致性蒸馏：结果

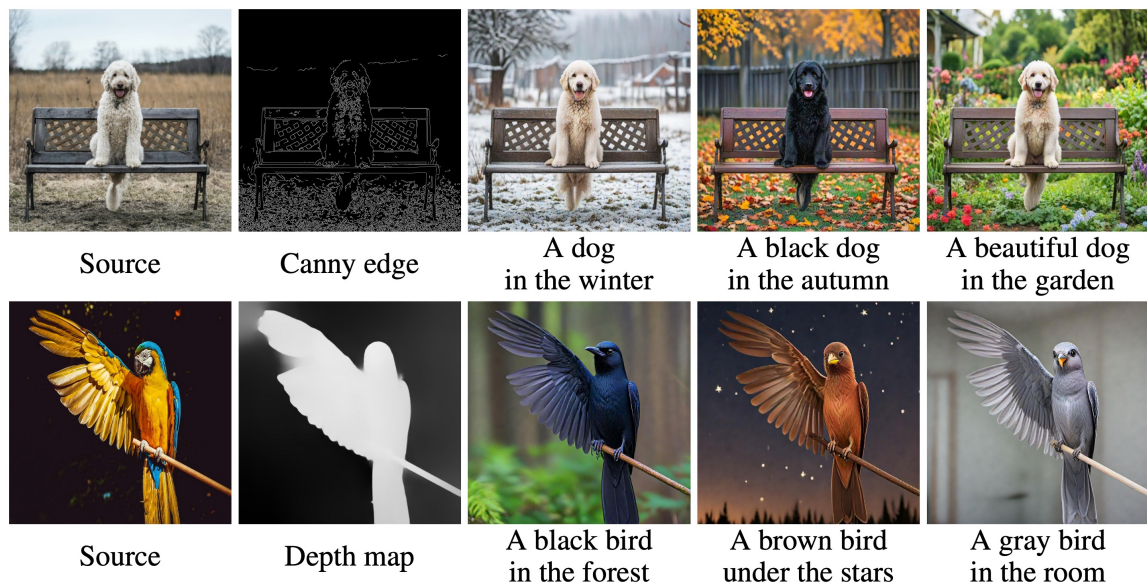


Figure 4: MLCM with ControlNet. Our MLCM can be incorporated into ControlNet pipeline and produce satisfactory results with 2 steps sampling.

## 与控制Net无缝结合：2步采样结果



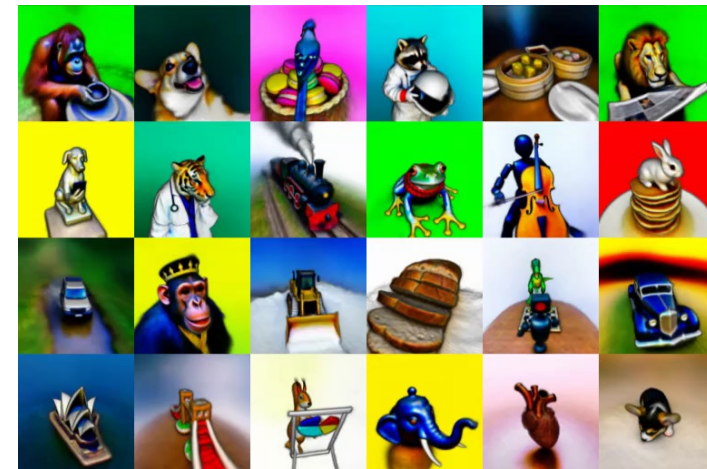
Figure 6: MLCM for Chinese-to-image generation. With 3 steps sampling, our MLCM model can produce images that align with Chinese semantic meaning. The first line presents images in general Chinese contexts, while the second line showcases images in specific Chinese cultural settings.

## 赋能中文-图像生成：3步采样结果

# 得分蒸馏 (score distillation) 和变分得分蒸馏 (variational score distillation, VSD)

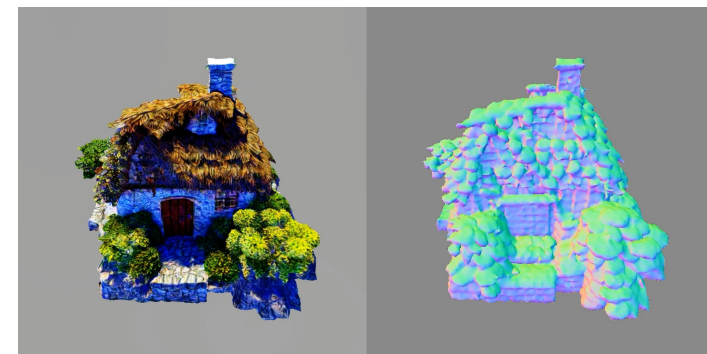
- 得分蒸馏: 基础扩散模型定义蒸馏的先验分布

$$\nabla_{\theta} \mathcal{L}_{\text{SDS}}(\theta) \triangleq \mathbb{E}_{t, \epsilon, c} \left[ \omega(t) (\epsilon_{\text{pretrain}}(\mathbf{x}_t, t, y) - \epsilon) \frac{\partial g(\theta, c)}{\partial \theta} \right]$$



- 变分得分蒸馏基于粒子的变分推断, 缓解过饱和问题

$$\frac{d\theta_{\tau}}{d\tau} = -\mathbb{E}_{t, \epsilon, c} \left[ \omega(t) \left( \underbrace{-\sigma_t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | y^c)}_{\text{score of noisy real images}} - \underbrace{(-\sigma_t \nabla_{\mathbf{x}_t} \log q_t^{\mu_{\tau}}(\mathbf{x}_t | c, y))}_{\text{score of noisy rendered images}} \right) \frac{\partial g(\theta_{\tau}, c)}{\partial \theta_{\tau}} \right]$$



已知基础扩散模型分布    未知渲染图像分布, 采用 LoRA 估计

问题: 收敛慢, 需要多阶段 (NeRF generation, geometry refinement, and texture refinement) 训练



## 变分得分蒸馏的实际改进：linearization+lookahead ( $L^2$ )

- 根据当前**3D**状态，**LoRA**模型向前看**1**步，保证对渲染数据分布的拟合

$$\nabla_{\theta_i} \mathcal{L}_{L-VSD}(\theta_i) = \mathbb{E}_{t, \epsilon, c} \left[ w(t) (\epsilon_{pretrain}(x_t, t, y) - \epsilon_{\phi_{i+1}}(x_t, t, c, y)) \frac{\partial g(\theta, c)}{\partial \theta} \Big|_{\theta=\theta_i} \right]$$

- 线性化**LoRA**模型防止过拟合

$$\epsilon_{\phi}^{\text{lin}}(x_t, t, c, y) = \epsilon_{\phi_i}(x_t, t, c, y) + (\phi - \phi_i) J_{\phi_i}^T(x_t, t, c, y)$$

- 基于**forward-mode autodiff**实现**VJP**

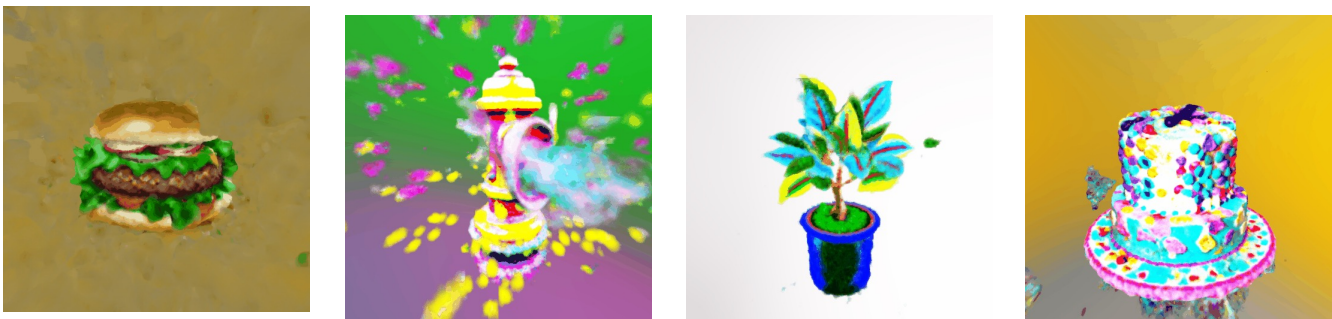


# L<sup>2</sup>-VSD: 结果

L<sup>2</sup>-VSD



ESD



VSD



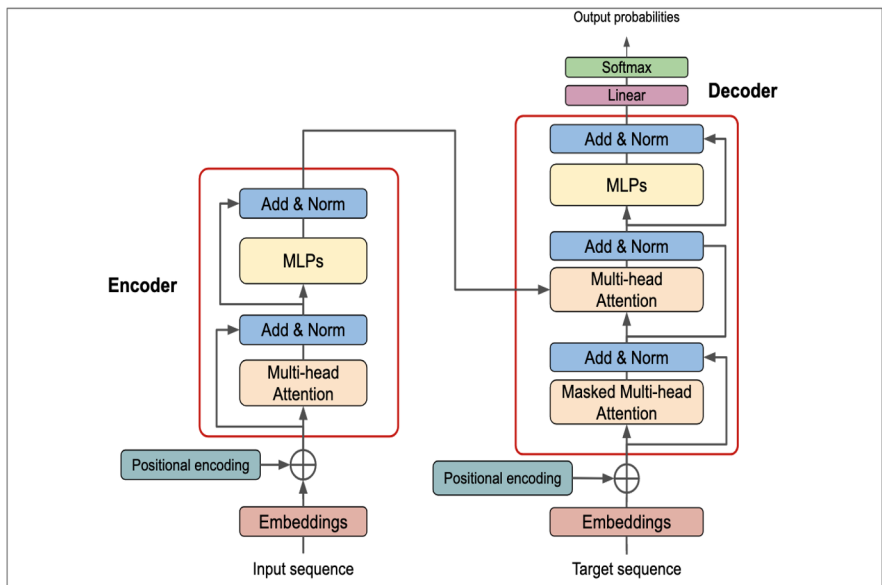
一个训练阶段即可生成高质量的结果



# AIGC大模型的高效推理方法

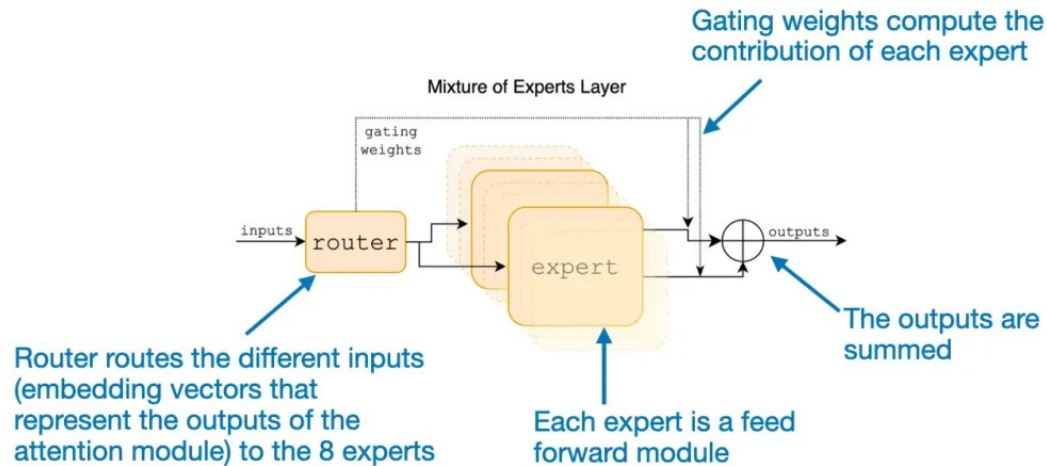
- 大语言模型的并行推理方法（15min）
  - 投机解码
  - 一致性蒸馏
- 大扩散模型的低步数推理方法（15min）
  - 采样器设计
  - 一致性蒸馏
- 大模型架构、序列状态、缓存等方面的优化方法（5min）

# 模型架构上的优化

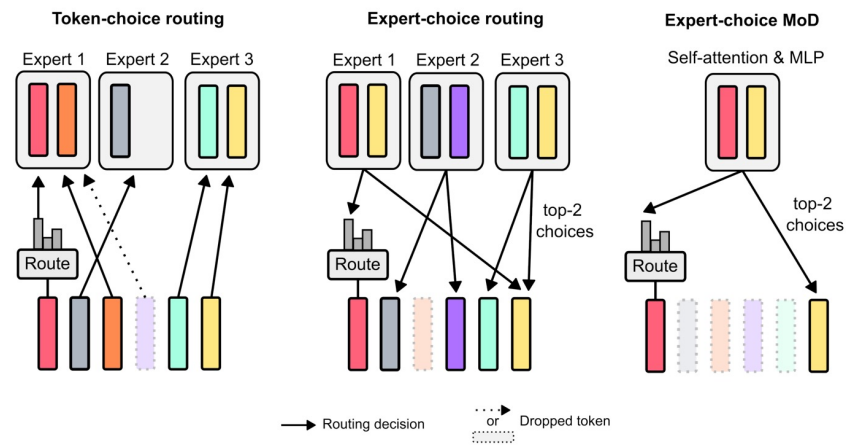


架构趋于统一: **Transformer**

问题: 如何进一步高效提升模型  
表达能力?



专家混合 (MoE): 稀疏激活专家, 保证 model capacity, 同时降低计算开销。代表: **Mixtral-8\*7B**



深度混合 (MoD): 从 token routing 到 expert routing, 保障负载平衡, 但难进行 casual modeling

# 模型架构上的优化

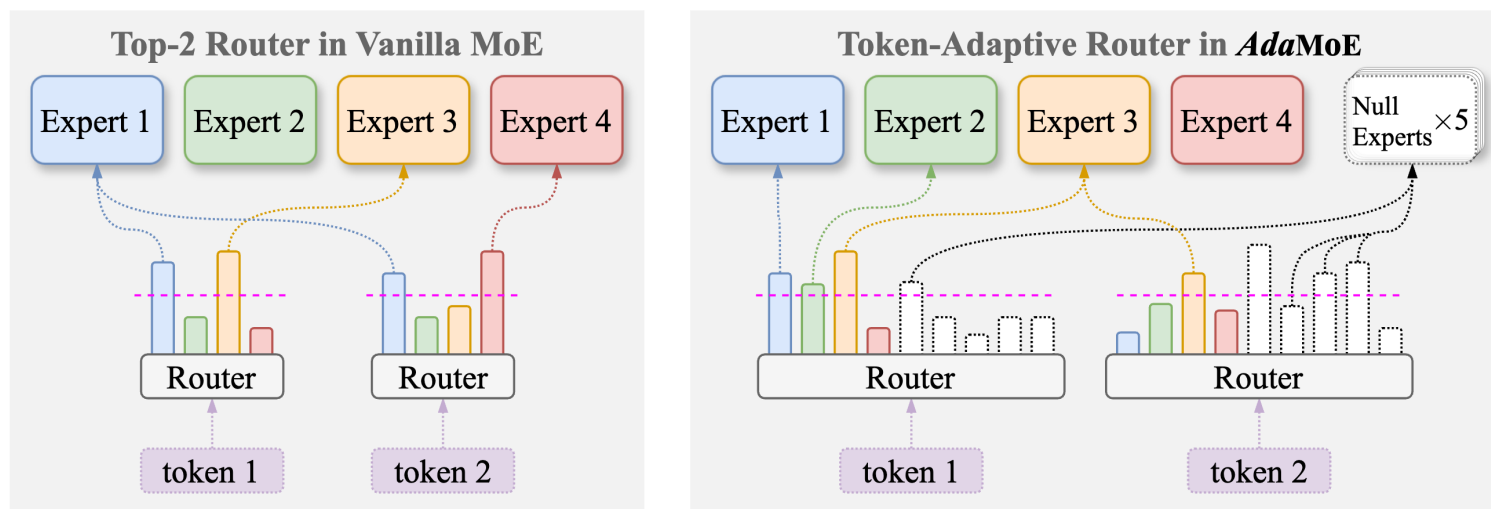


Figure 2: Comparison of Routing Mechanisms: Vanilla MoE vs. *AdaMoE*. **Left:** In the Vanilla MoE model, each token selects the top-2 experts based on the highest router probabilities. **Right:** The *AdaMoE* model introduces an additional set of *null experts*. Each token selects the top-4 experts, which can include both true and null experts. For example, token 1 selects three true experts, while token 2 selects only one true expert. Despite this variation, the average number of true experts selected per token remains two, maintaining parity with the vanilla method.

**AdaMoE:** 引入空专家，增加灵活度，不同token可以选0-4个真专家，负载均衡损失保证空专家使用率，建模灵活性优于MoD



# 模型架构上的优化

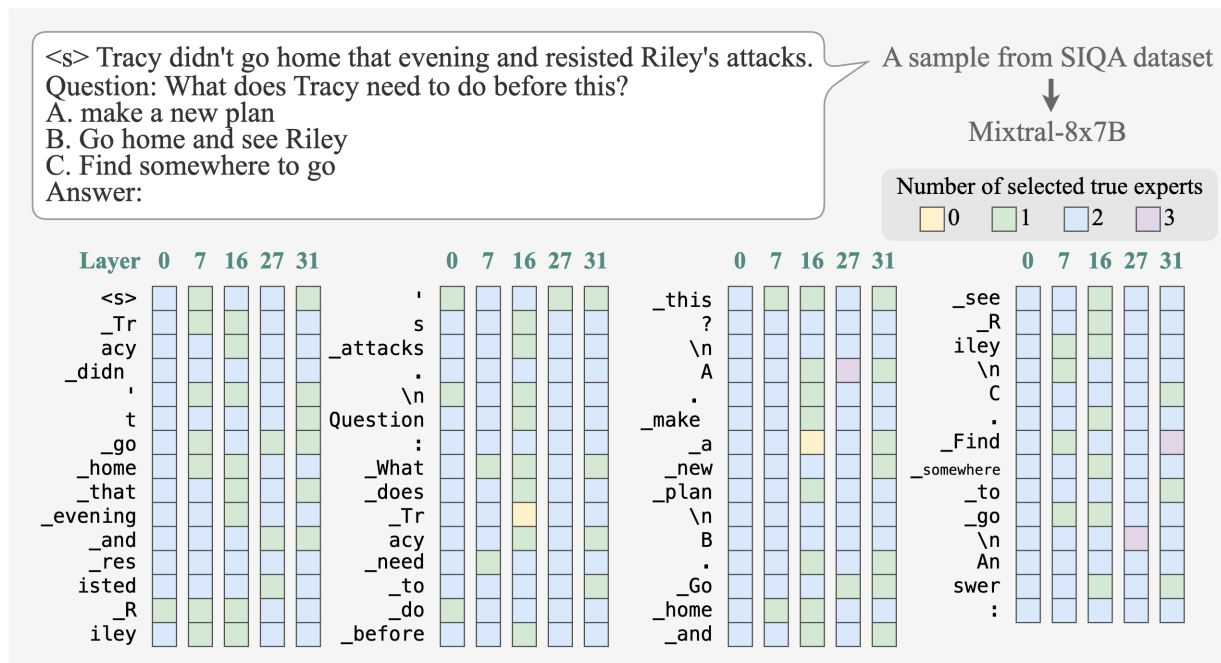
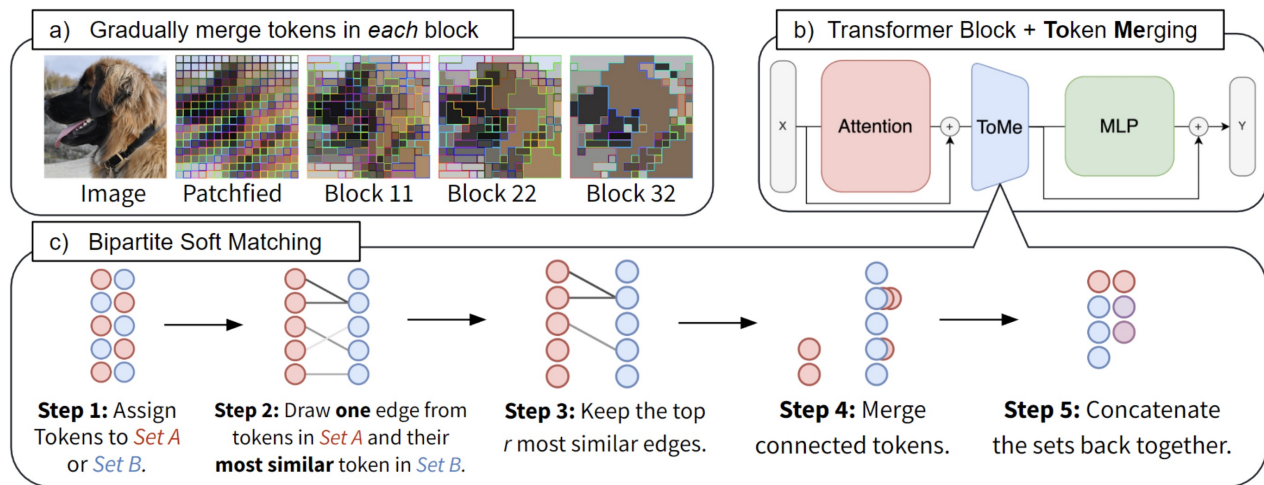


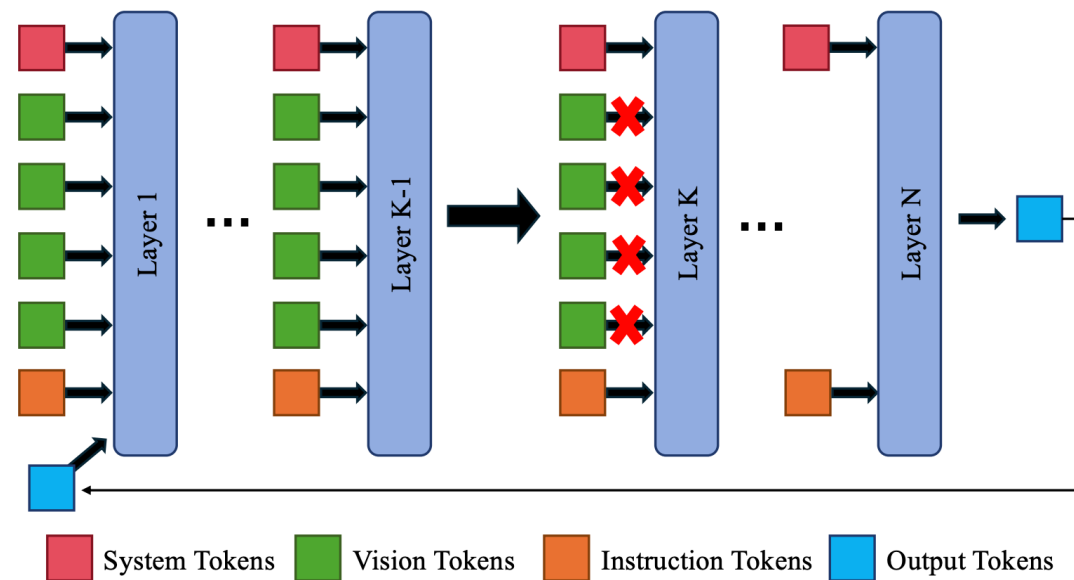
Figure 1: Inference on a sample from SIQA dataset after fine-tuning Mixtral-8x7b with *AdaMOE*. The image shows the number of true experts selected for each token in different layers.

对**Mixtral-8\*7B**简单微调，即可收获大量选择**1**个真专家的**token**，性能不减的情况下降低至多**20%**的**FLOPs**

# Token序列上的优化



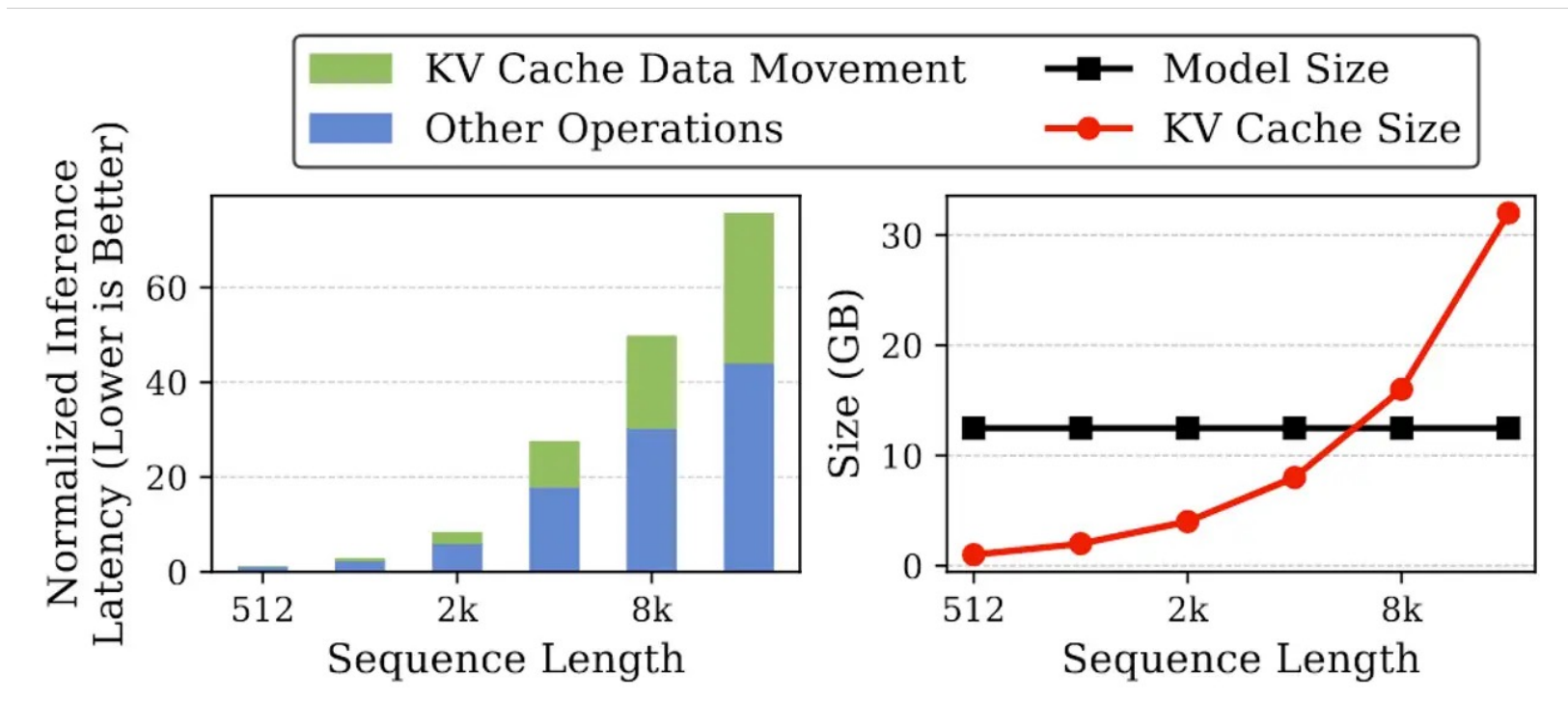
基于相似度对token序列进行合并，缩减序列长度，降低计算开销，可用于图像和视频生成



**Multi-modal**大模型中，中间视觉tokens在深层中获得的注意力极少，其信息转移到文本tokens中，因此若干层（如：16）之后，可直接丢弃视觉tokens

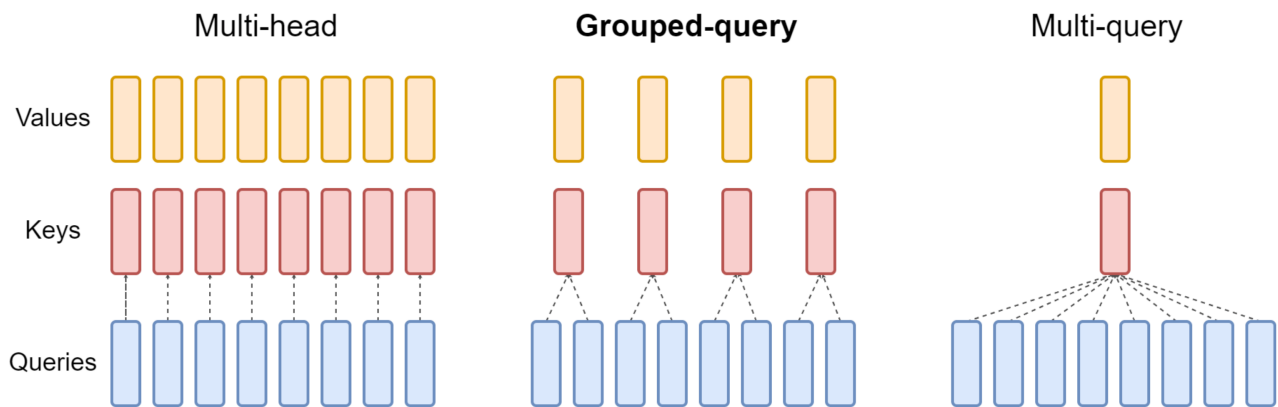


# 缓存 (KV Cache) 上的优化



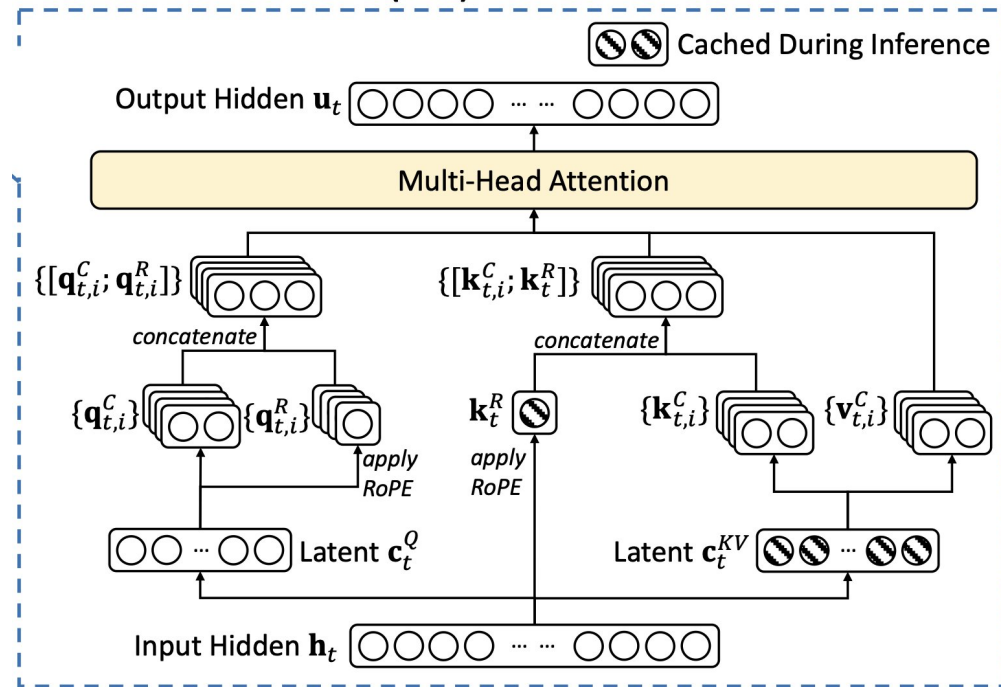
长序列场景下，访存 (KV Cache) 成为瓶颈 (MPT-7B)

# 缓存 (KV Cache) 上的优化



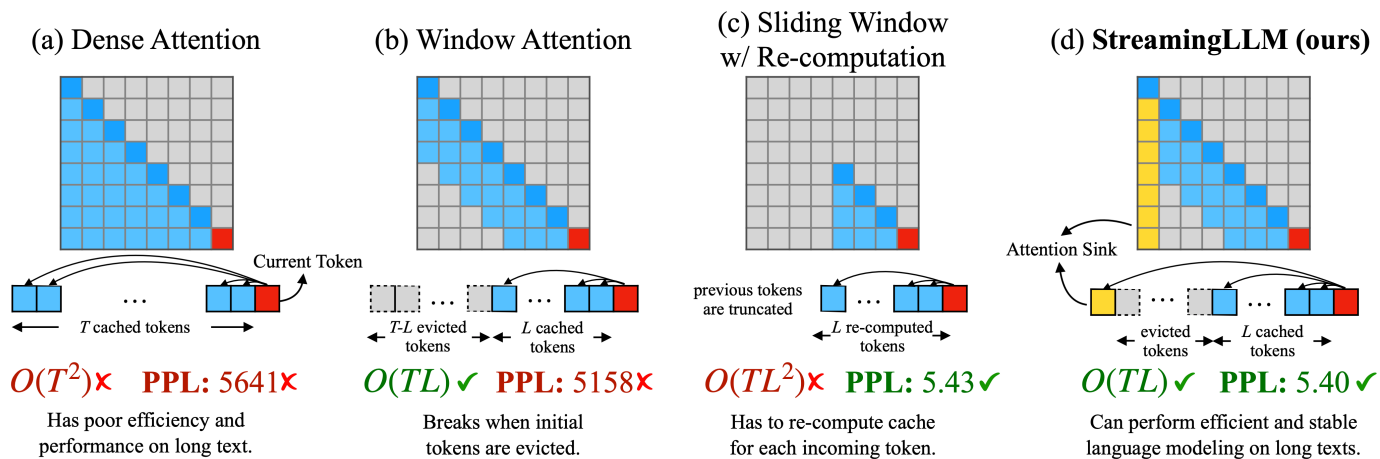
**Grouped-Query Attention:** 将查询向量分组处理，共享KV

## Multi-Head Latent Attention (MLA)



**Multi-head Latent Attention:** 将KV cache投影到低维隐向量，用于DeepSeek v2

# 缓存 (KV Cache) 上的优化



**StreamingLLM: window attention+起始位置 attention, 弥补的window attention精度损失**

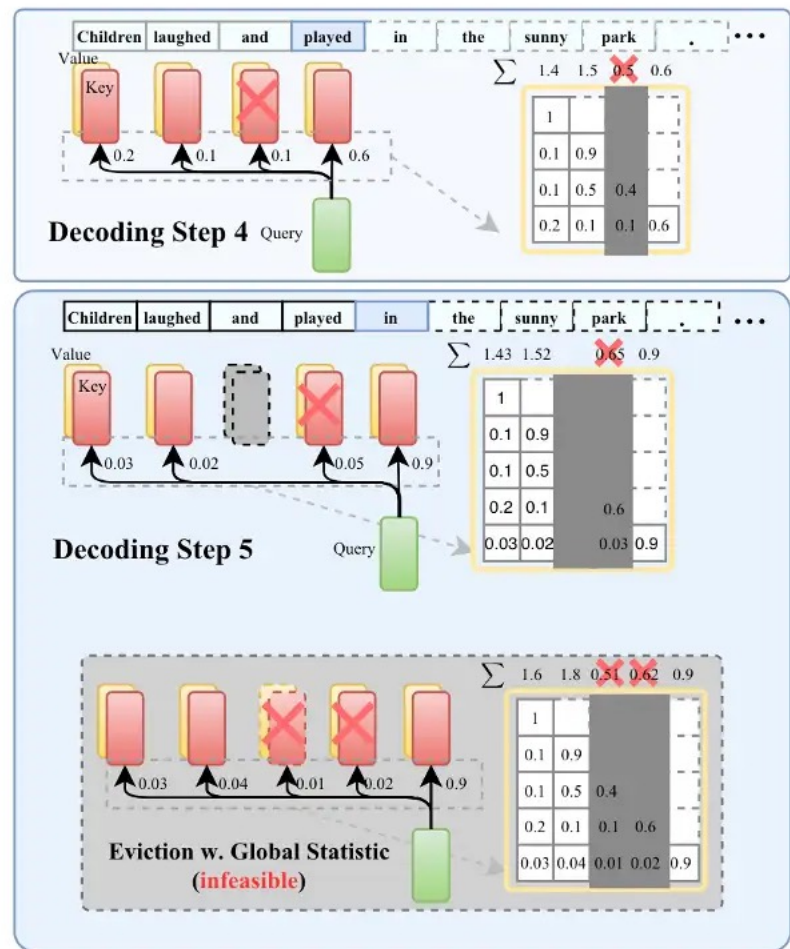


Figure 3: Illustration of Algorithm 1 during two consecutive decoding steps.

**H2O: 贪心地丢掉当前最不被注意的token**





# AIGC大模型的高效推理方法：未来研究方向

- 基础：
  - 模型架构、学习方法、采样算法
  
- 场景：
  - 多模态大模型
  - 视频生成模型（世界模型）

感谢各位专家！ 敬请批评指正！

邮箱：[zhijied@sjtu.edu.cn](mailto:zhijied@sjtu.edu.cn)

主页：<https://thudzj.github.io/>

