

Text Diffusion Models, Multimodal LLMs

COMP7607 — Lecture 9

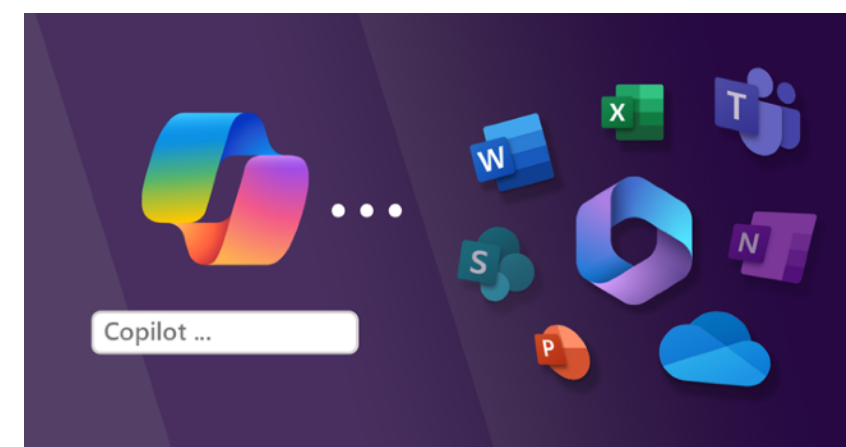
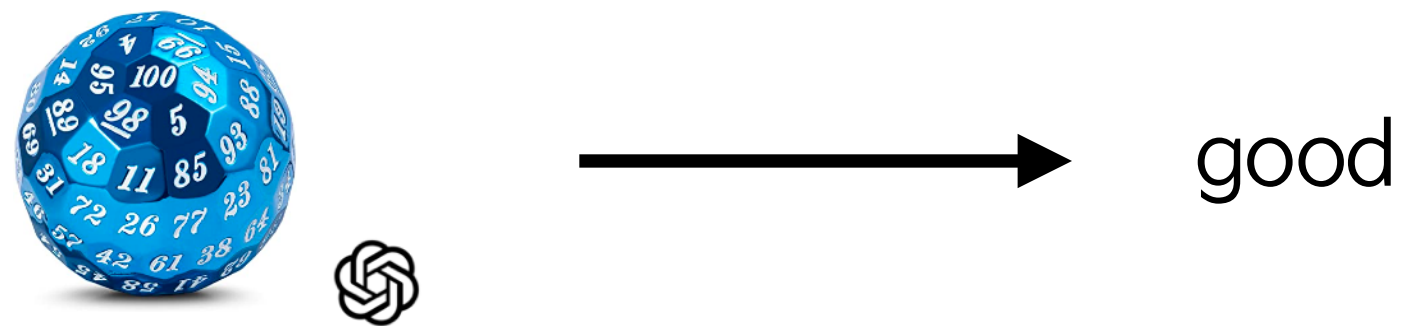
Lingpeng Kong

Department of Computer Science, The University of Hong Kong

Generate a sentence

$$p(\boldsymbol{x}) = \prod_i p(x_i | \boldsymbol{x}_{<i})$$

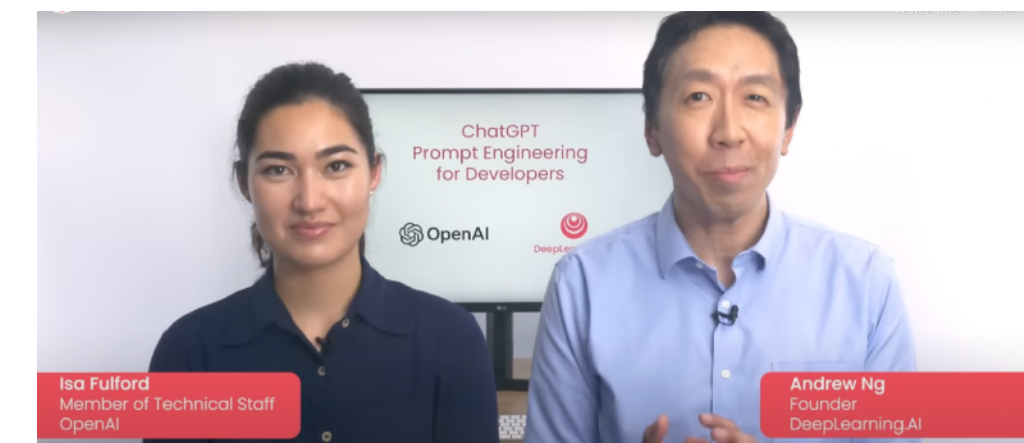
Don't just believe that it is because something is trendy that it is ...



LLMs as Agents



LLM Reasoning



Prompt Engineering

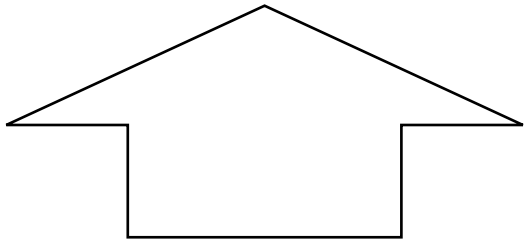
Non-autogressive Text Generation

Tokyo is the largest city in the world



Iterative Refinement

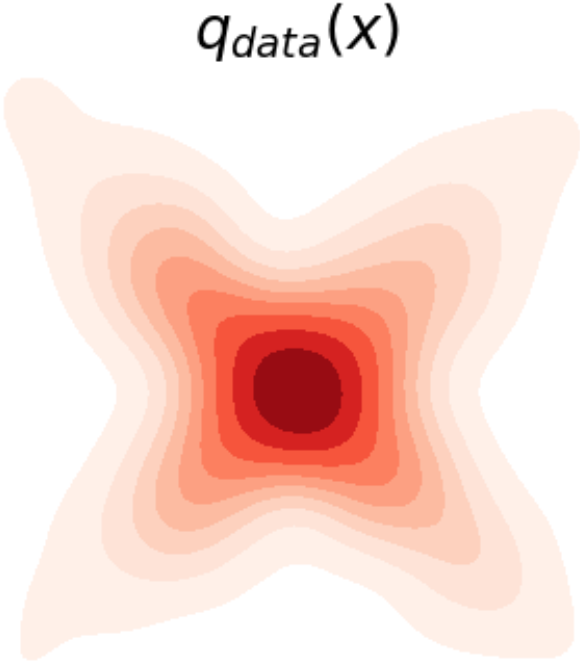
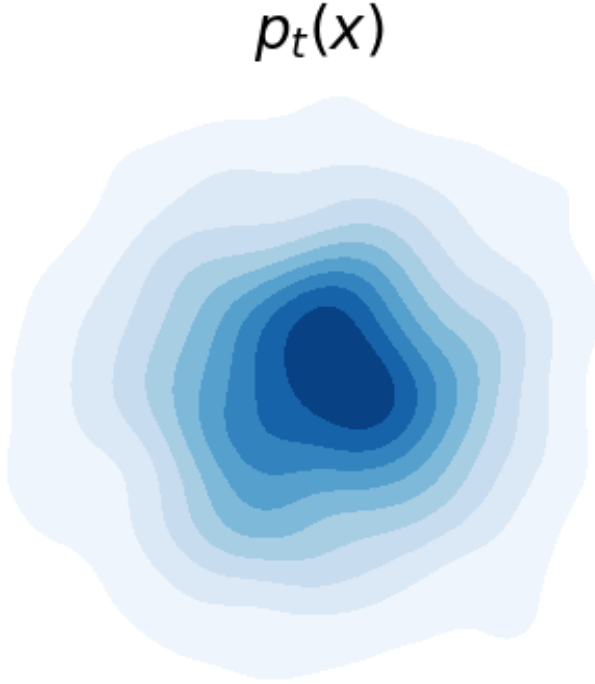
Shanghai is the largest city in the world



Parallel Generation

4+10 14 14-12 2 13*2=24
3*13-39,9+10-19,39-19-24
1+10-11 13-11=2,7*2=24

[mask] [mask] [mask] [mask] [mask] [mask] [mask] [mask]



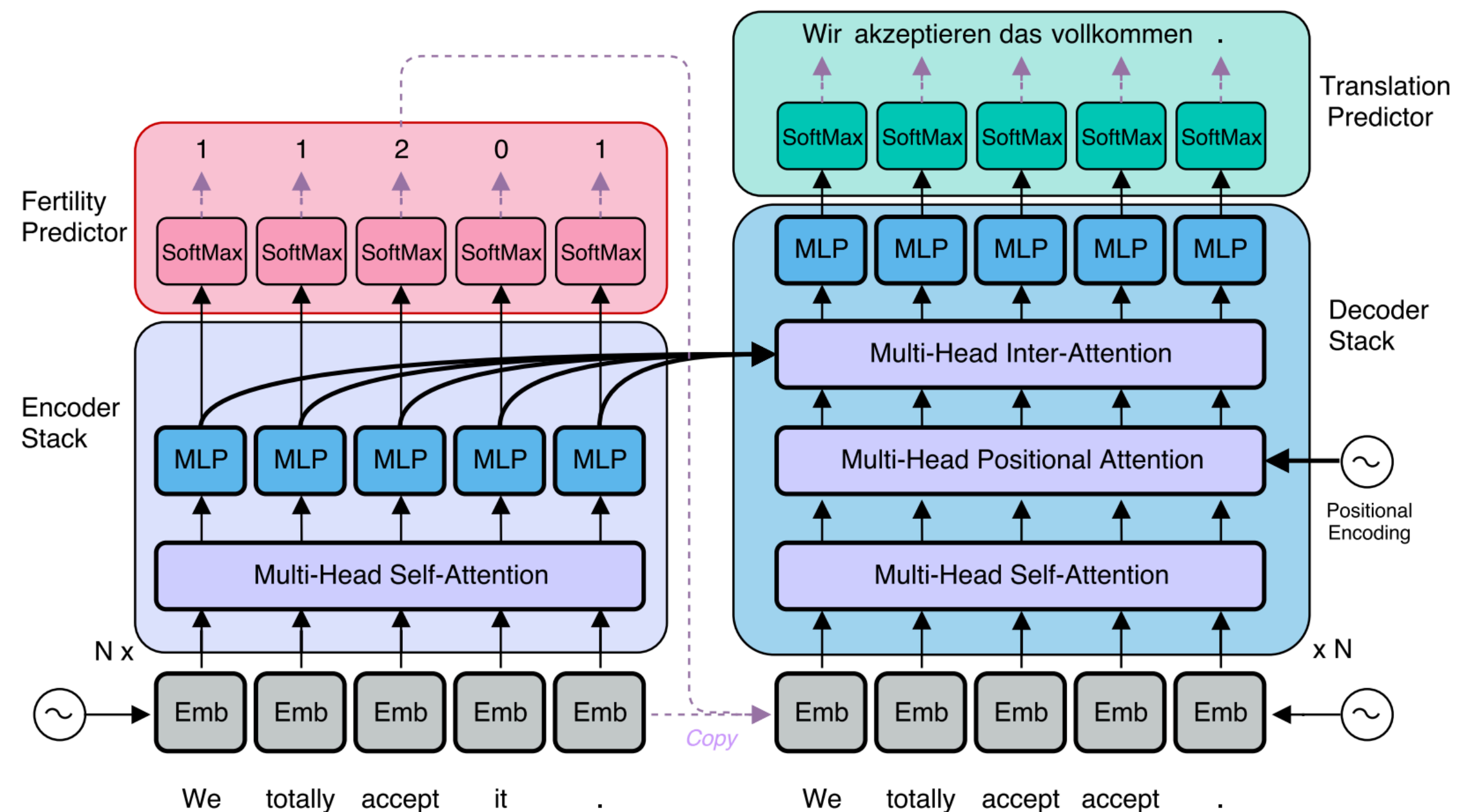
In a different “family” of $p(x)$ than the autoregressive ones

Non-autogressive Neural Machine Translation

Models	WMT14	
	En→De	De→En
NAT	17.35	20.62
NAT (+FT)	17.69	21.47
NAT (+FT + NPD $s = 10$)	18.66	22.41
NAT (+FT + NPD $s = 100$)	19.17	23.20
Autoregressive ($b = 1$)	22.71	26.39
Autoregressive ($b = 4$)	23.45	27.02

many thanks
thank you → thank thanks

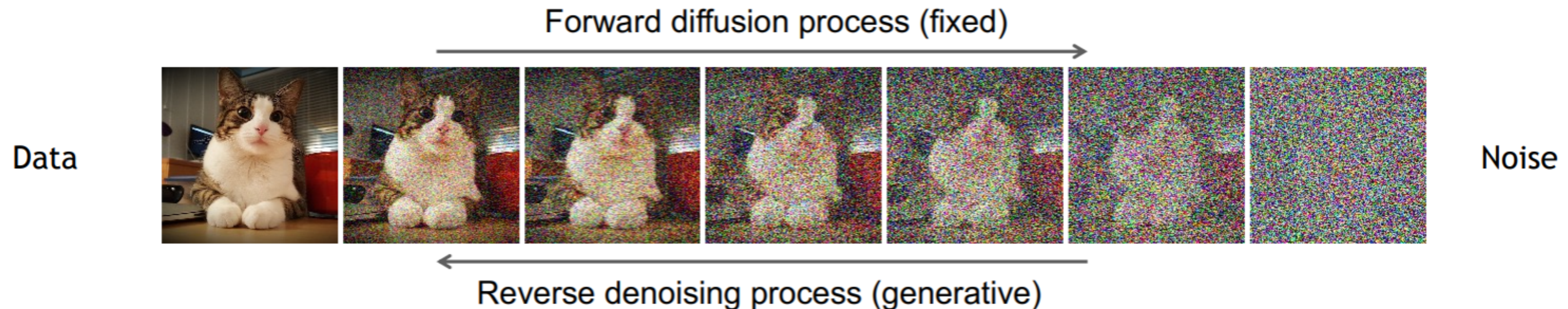
“Modality Conflicts”



Diffusion Models

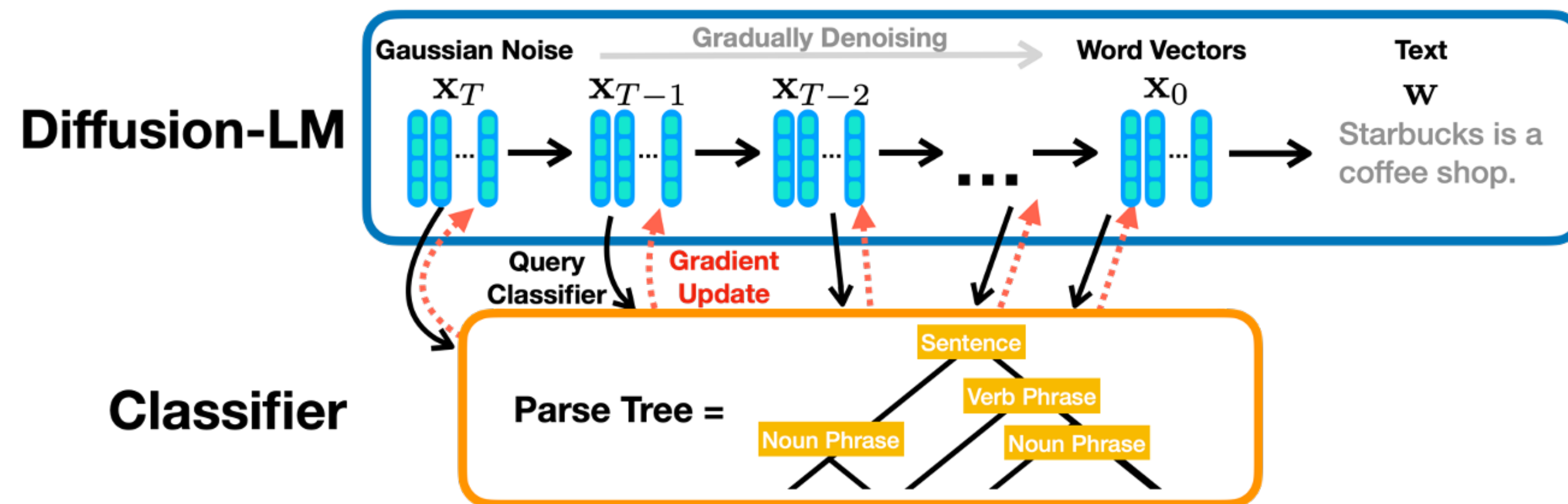
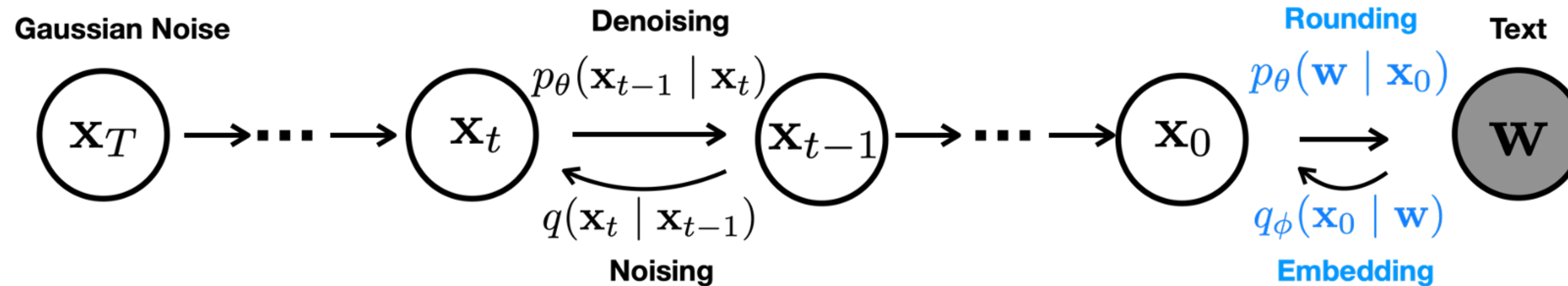
Forward-backward formulation

The forward process gradually injects noise to the input
The backward process denoises to recover the original data

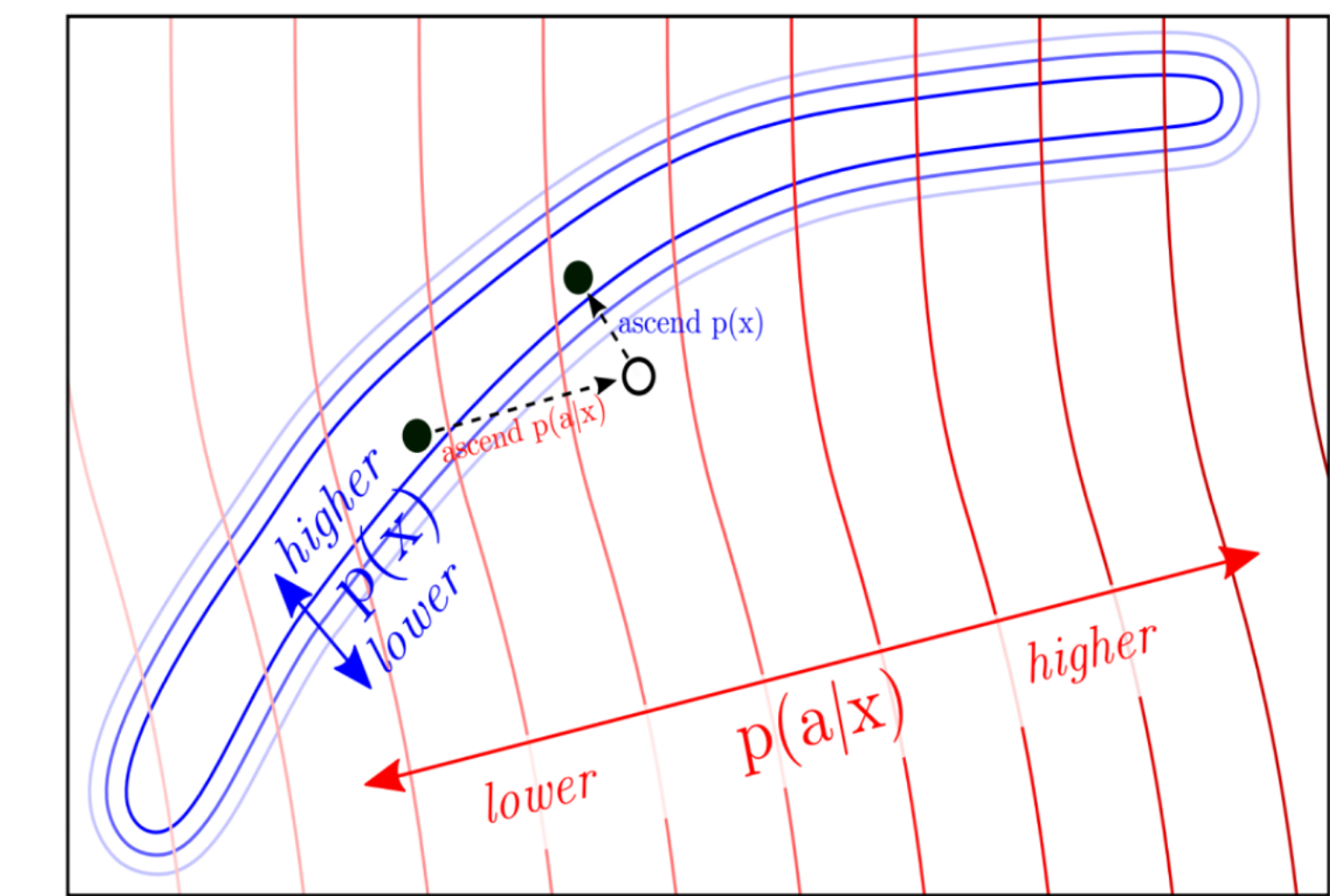
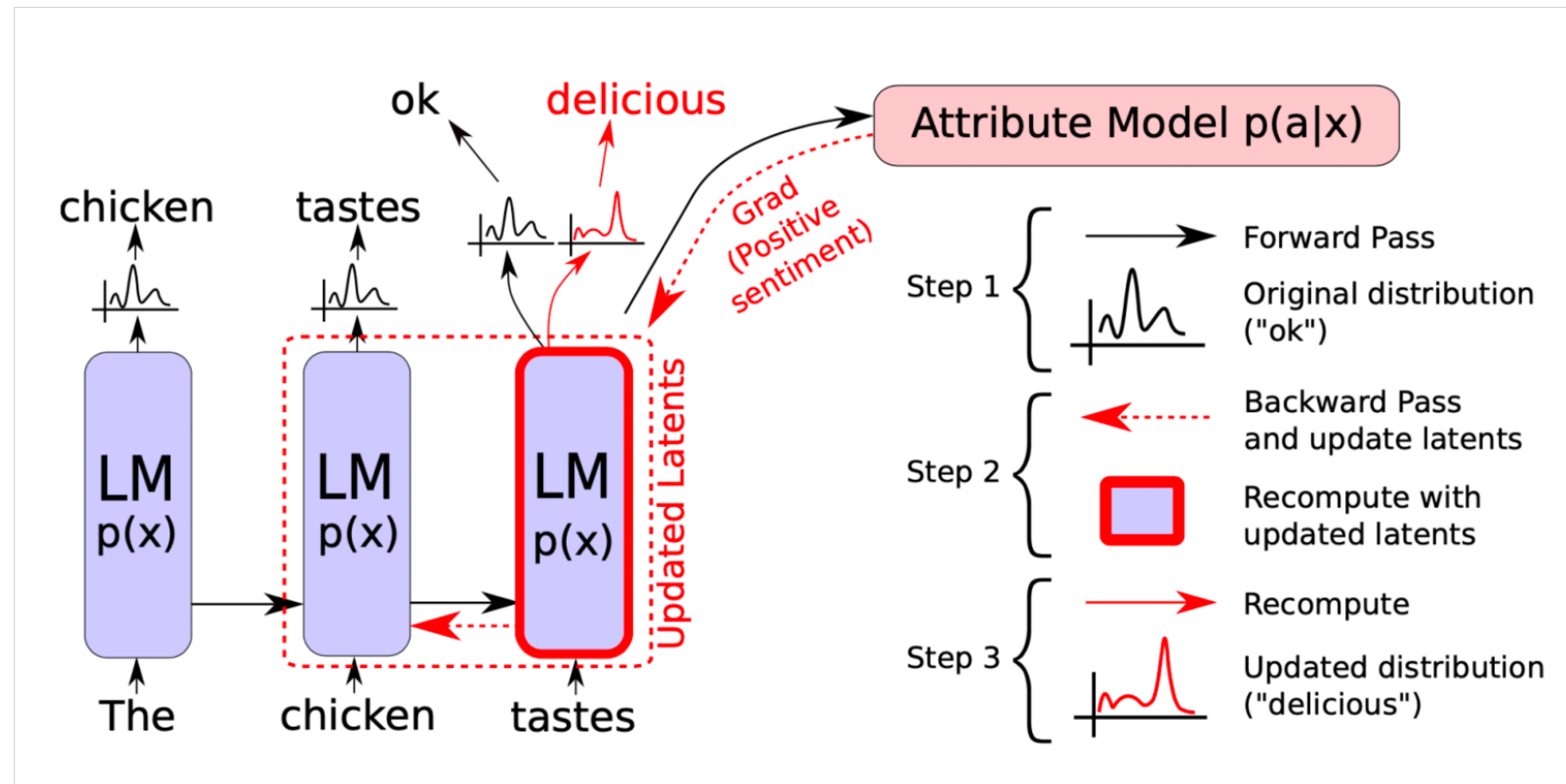


“Creating noise from data is easy; creating data from noise is generative modeling.” (Song et al., 2021)

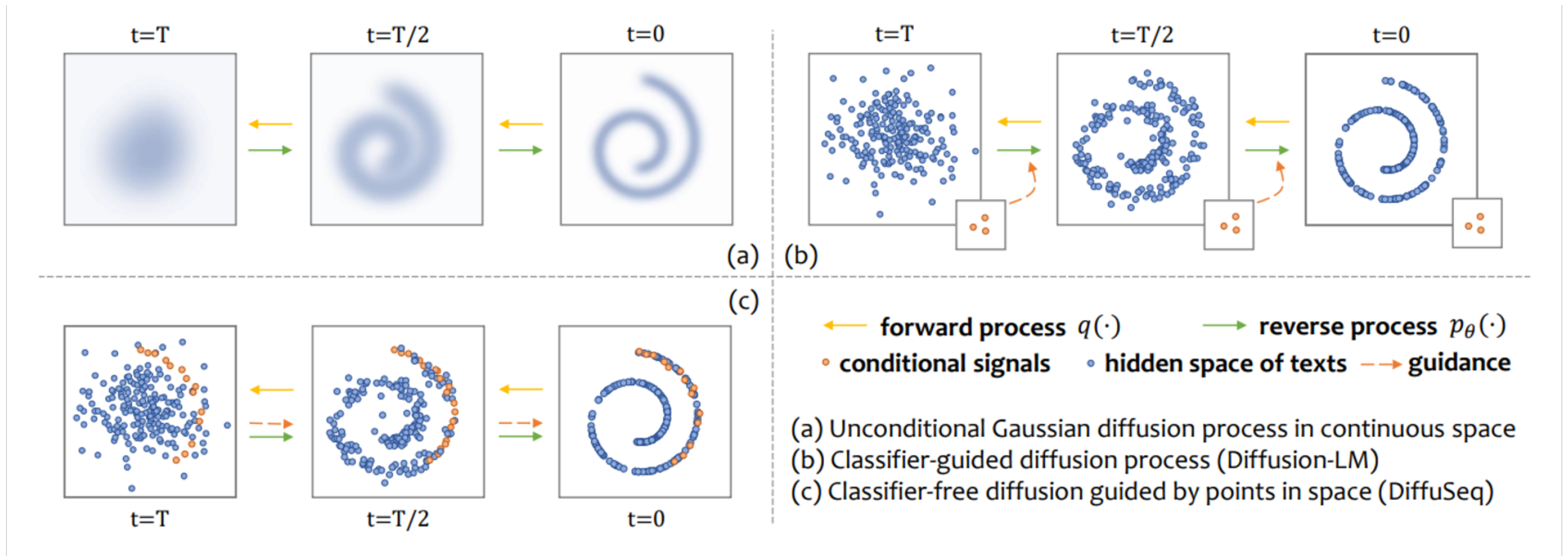
Text Diffusion Models



Classifier Guidance

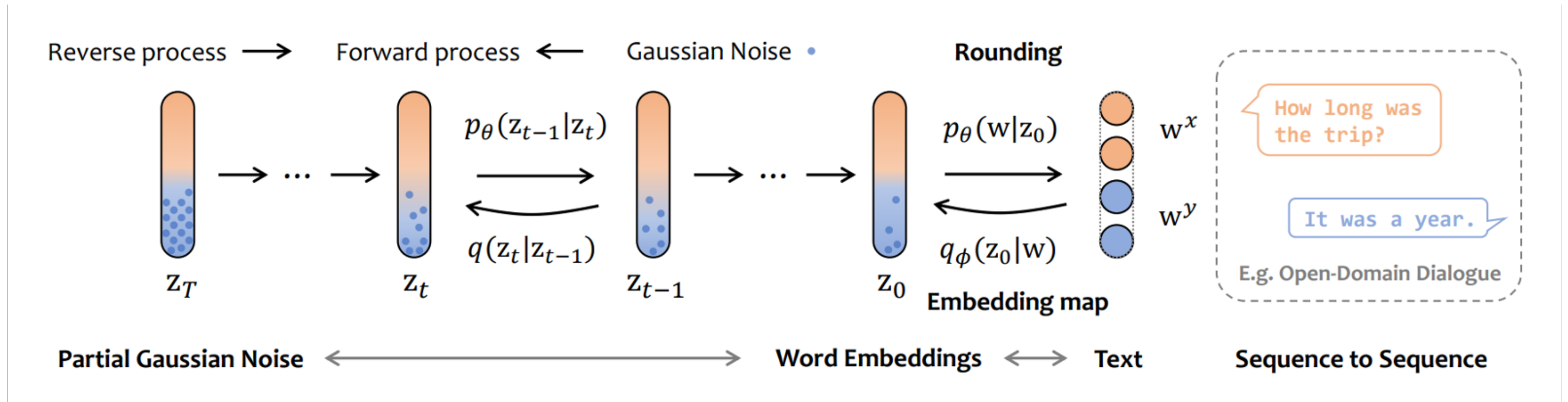


Sequence to Sequence Text Diffusion




“Seq2Seq” tasks: $\mathbf{x} \rightarrow \mathbf{y}$

Sequence to Sequence Text Diffusion



$$\mathcal{L}_{\text{VLB}} = \mathbb{E}_{q(\mathbf{z}_{1:T}|\mathbf{z}_0)} \left[\underbrace{\log \frac{q(\mathbf{z}_T|\mathbf{z}_0)}{p_\theta(\mathbf{z}_T)}}_{\mathcal{L}_T} + \sum_{t=2}^T \underbrace{\log \frac{q(\mathbf{z}_{t-1}|\mathbf{z}_0, \mathbf{z}_t)}{p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t)}}_{\mathcal{L}_{t-1}} + \underbrace{\log \frac{q_\phi(\mathbf{z}_0|\mathbf{w}^{x \oplus y})}{p_\theta(\mathbf{z}_0|\mathbf{z}_1)}}_{\mathcal{L}_0} - \underbrace{\log p_\theta(\mathbf{w}^{x \oplus y}|\mathbf{z}_0)}_{\mathcal{L}_{\text{round}}} \right].$$

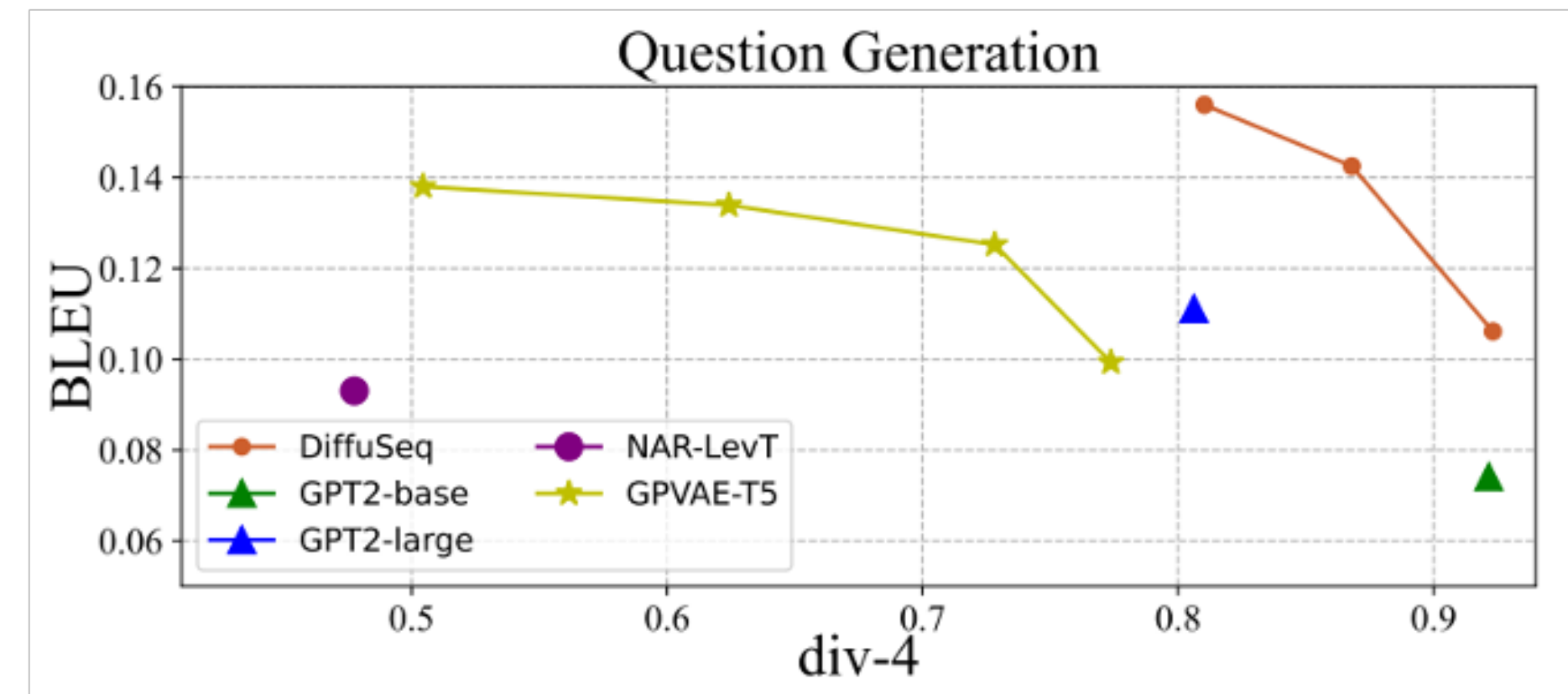
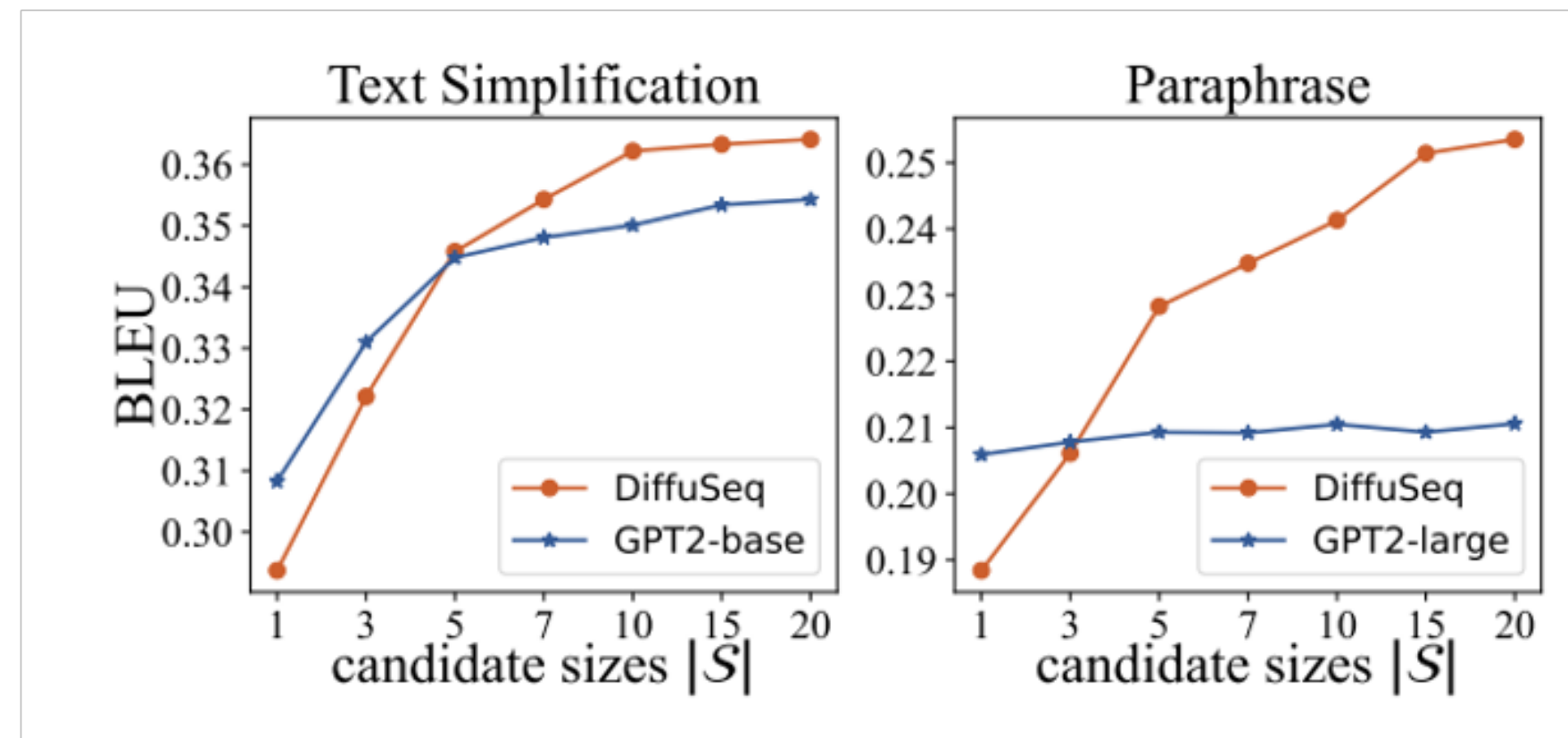
\times



$$\min_{\theta} \mathcal{L}_{\text{VLB}} = \min_{\theta} \left[\sum_{t=2}^T ||\mathbf{y}_0 - \tilde{f}_\theta(\mathbf{z}_t, t)||^2 + ||\text{EMB}(\mathbf{w}^y) - \tilde{f}_\theta(\mathbf{z}_1, 1)||^2 + \mathcal{R}(||\mathbf{z}_0||^2) \right]$$

Sequence to Sequence Text Diffusion

Diversity Ensures Quality



Statement: *The Japanese yen is the official and only currency recognized in Japan.*

Question: *What is the Japanese currency?*

GPVAE-T5

- * What is the japanese currency
- * What is the japanese currency
- * What is the japanese currency

NAR-LevT

- * What is the basic unit of currency for Japan ?
- * What is the basic unit of currency for Japan ?
- * What is the basic unit of currency for Japan ?

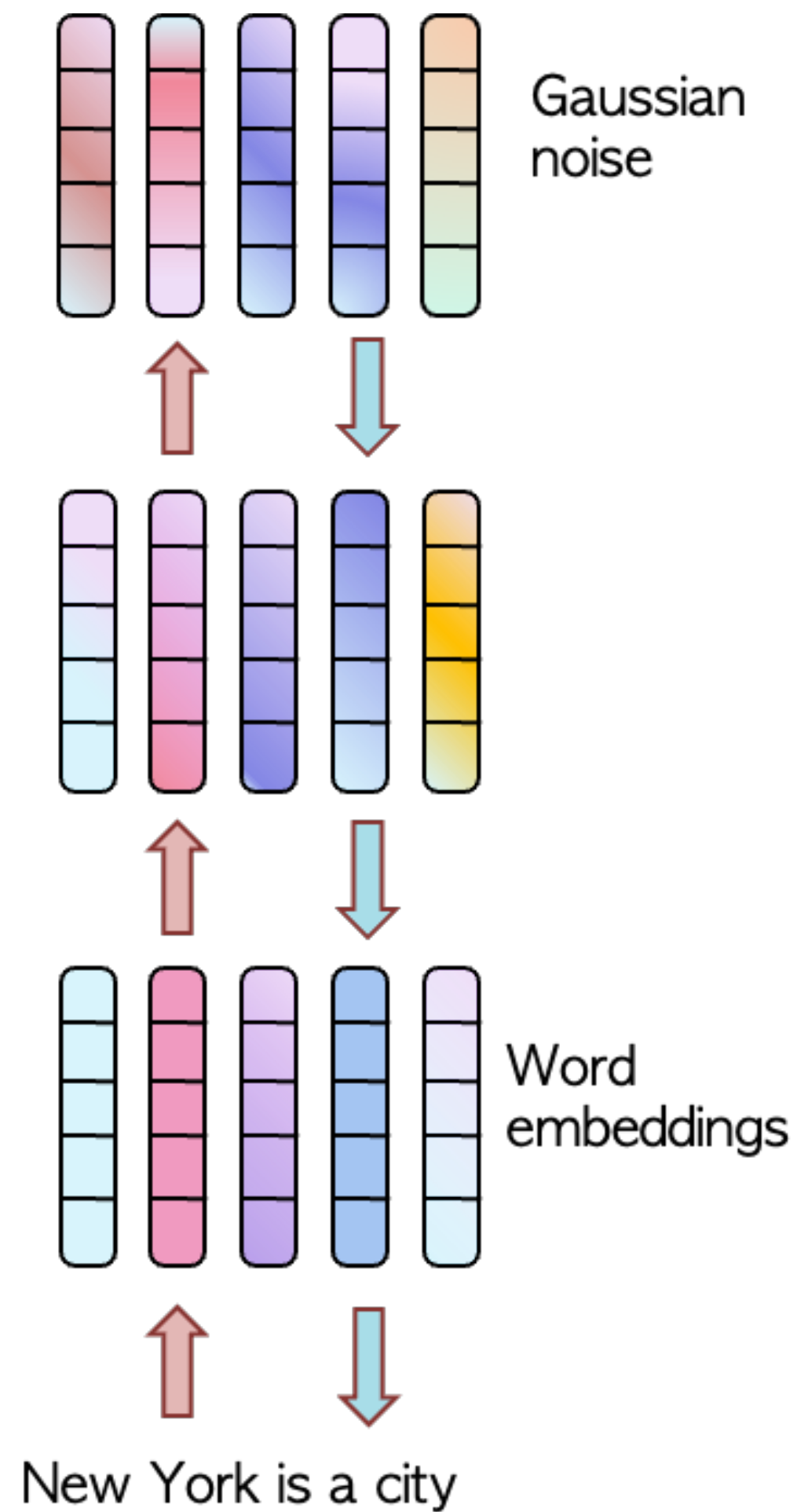
GPT2-large finetune

- * What is the basic unit of currency for Japan?
- * What is the Japanese currency
- * What is the basic unit of currency for Japan?

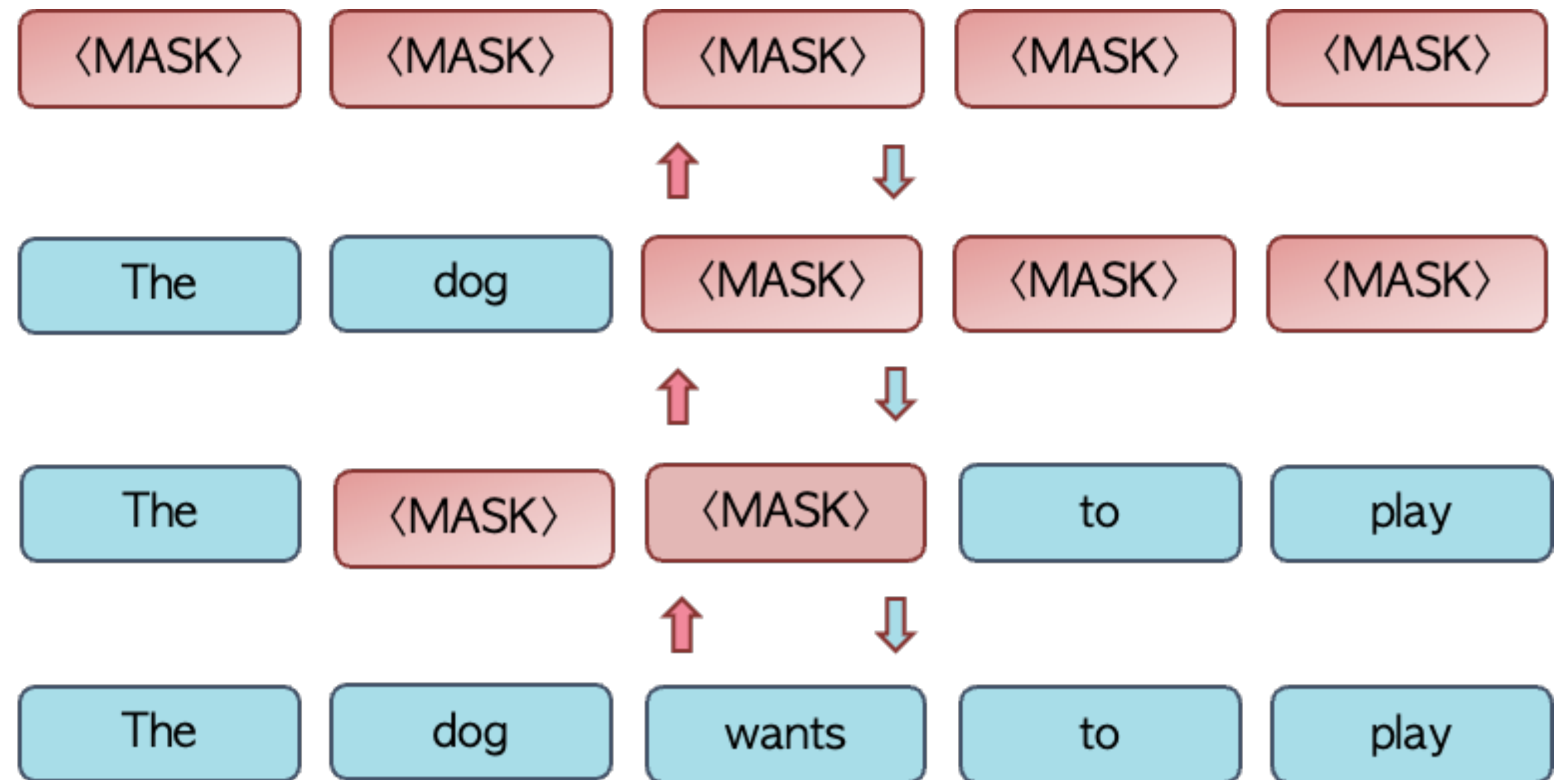
DiffuSeq

- * What is the Japanese currency
- * Which country uses the “yen yen” in currency
- * What is the basic unit of currency?

Text Diffusion Models



Continuous Diffusion



Discrete Diffusion

Discrete Diffusion - Forward Process

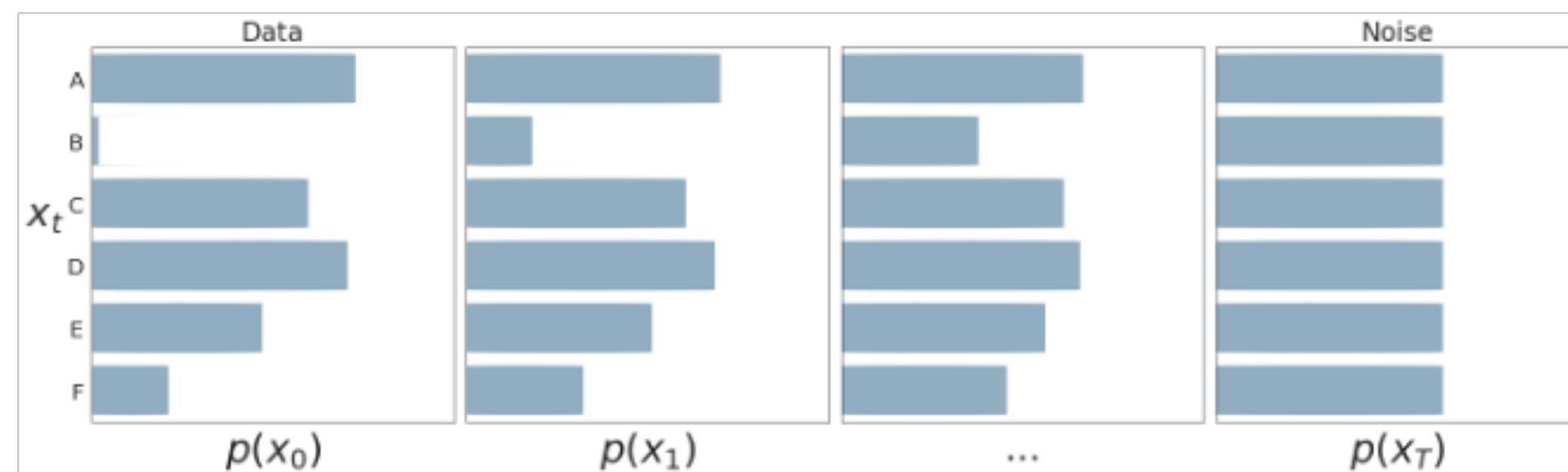
$$q(x_0) = p_{\text{data}}(x)$$

Noise: $q(x_t \mid x_{t-1}) = \beta_t x_{t-1} + (1 - \beta_t) q_{\text{noise}}$

Enjoys closed-form “marginal” distribution:

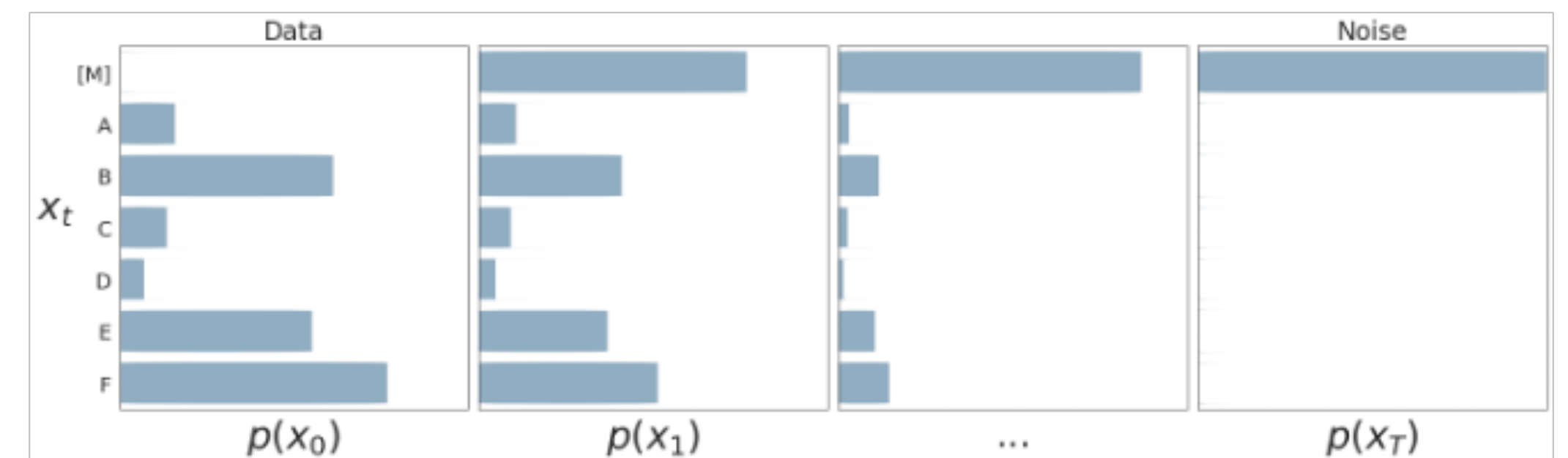
$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \alpha_t \mathbf{x}_0 + (1 - \alpha_t) q_{\text{noise}}, \quad \alpha_t := \prod_{i=1}^t \beta_i$$

“Multinomial”



$$q_{\text{noise}} = \frac{1}{K}$$

“Absorbing”



$$q_{\text{noise}} = \delta([M])$$

Discrete Diffusion - Learning

Goal: $p_{\theta}(x_{t-1} \mid x_t) \approx q(x_{t-1} \mid x_t)$

$$\begin{aligned} \log p_{\theta}(\mathbf{x}_0) &\geq \mathbb{E}_{q(\mathbf{x}_1, \dots, \mathbf{x}_T \mid \mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T)}{q(\mathbf{x}_1, \dots, \mathbf{x}_T \mid \mathbf{x}_0)} \right] \\ &= \mathbb{E}_q \left[\underbrace{\log p_{\theta}(\mathbf{x}_0 \mid \mathbf{x}_1)}_{\mathcal{L}_1(\theta)} - \sum_{t=2}^T \underbrace{\text{KL}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t))}_{\mathcal{L}_t(\theta)} \right] + \text{const.} \end{aligned}$$

Enjoys closed-form “marginal” distribution:

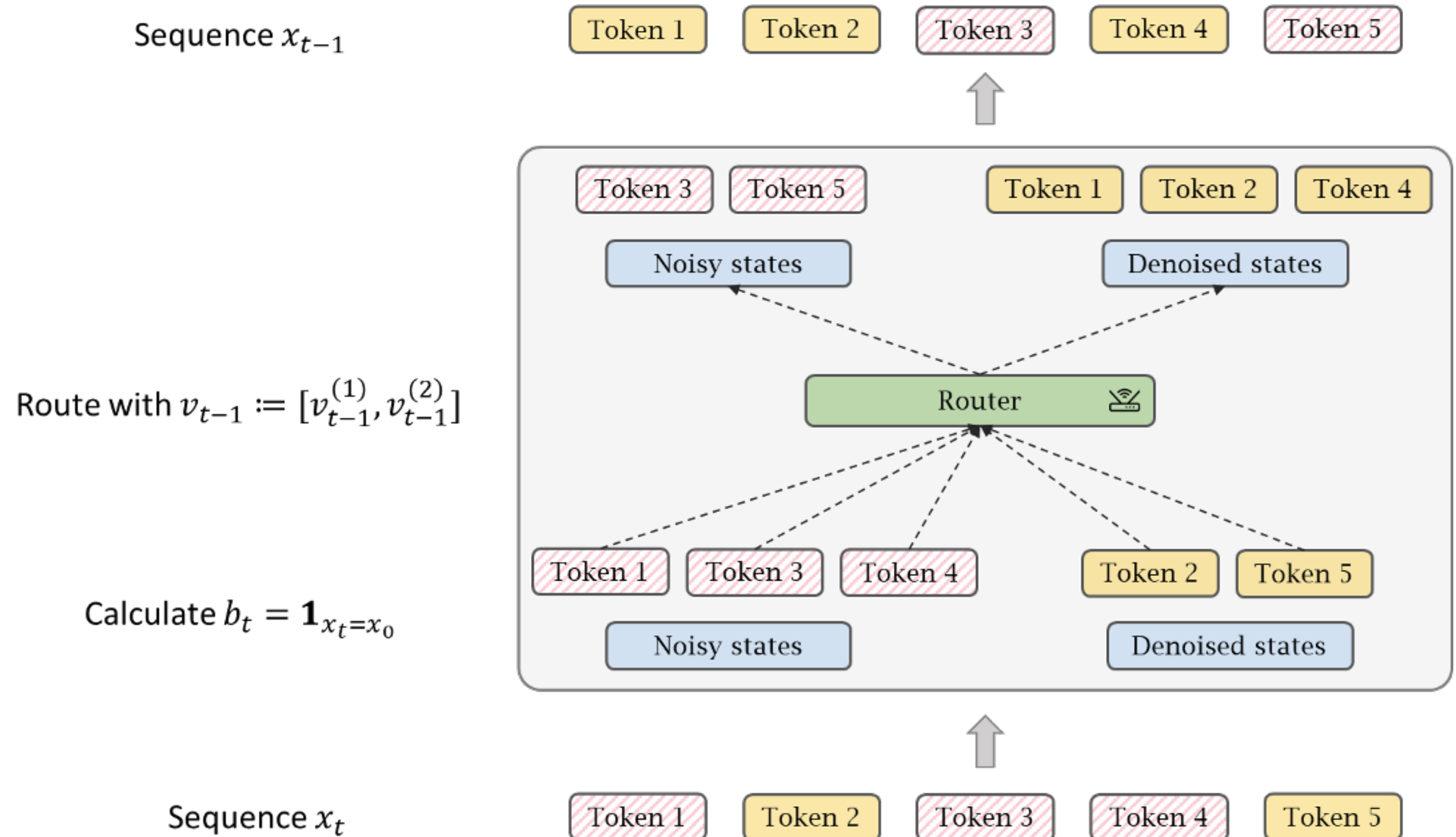
“Multinomial” $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \frac{(\beta_t \mathbf{x}_t + (1 - \beta_t) \frac{\mathbf{1}}{K}) \odot (\alpha_{t-1} \mathbf{x}_0 + (1 - \alpha_{t-1}) \frac{\mathbf{1}}{K})}{\mathbf{x}_t^{\top} (\alpha_t \mathbf{x}_0 + (1 - \alpha_t) \frac{\mathbf{1}}{K})}$

“Absorbing” $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \begin{cases} \mathbf{x}_t, & \text{if } \mathbf{x}_t \neq \mathbf{e}_{[M]}, \\ \frac{1 - \alpha_{t-1}}{1 - \alpha_t} \mathbf{e}_{[M]} + \frac{\alpha_{t-1} - \alpha_t}{1 - \alpha_t} \mathbf{x}_0, & \text{if } \mathbf{x}_t = \mathbf{e}_{[M]} \end{cases}$

Sampling Reparameterization

Reveals a latent routing mechanism

$$\begin{aligned}
 b_t &= \mathbf{1}_{\mathbf{x}_t = \mathbf{x}_0} \\
 v_{t-1}^{(1)} &\sim \text{Bernoulli}(\lambda_{t-1}^{(1)}), \quad \mathbf{u}_t^{(1)} \sim q_{\text{noise}} \\
 v_{t-1}^{(2)} &\sim \text{Bernoulli}(\lambda_{t-1}^{(2)}), \quad \mathbf{u}_t^{(2)} \sim q_{\text{noise}}(\mathbf{x}_t) \\
 \mathbf{x}_{t-1} &= b_t \left[v_{t-1}^{(1)} \mathbf{x}_t + \left(1 - v_{t-1}^{(1)}\right) \mathbf{u}_t^{(1)} \right] + \\
 &\quad (1 - b_t) \left[v_{t-1}^{(2)} \mathbf{x}_0 + \left(1 - v_{t-1}^{(2)}\right) \mathbf{u}_t^{(2)} \right]
 \end{aligned}$$



Training RDMs

Trained with ELBO:

$$\mathcal{L}_t := -\lambda_{t-1}^{(2)} \sum_{n=1}^N \mathbf{1}_{\mathbf{x}_{t,n} \neq \mathbf{x}_{0,n}} \mathbf{x}_{0,n}^\top \log f(\mathbf{x}_{t,n}; \boldsymbol{\theta})$$

Time-dependent weighting

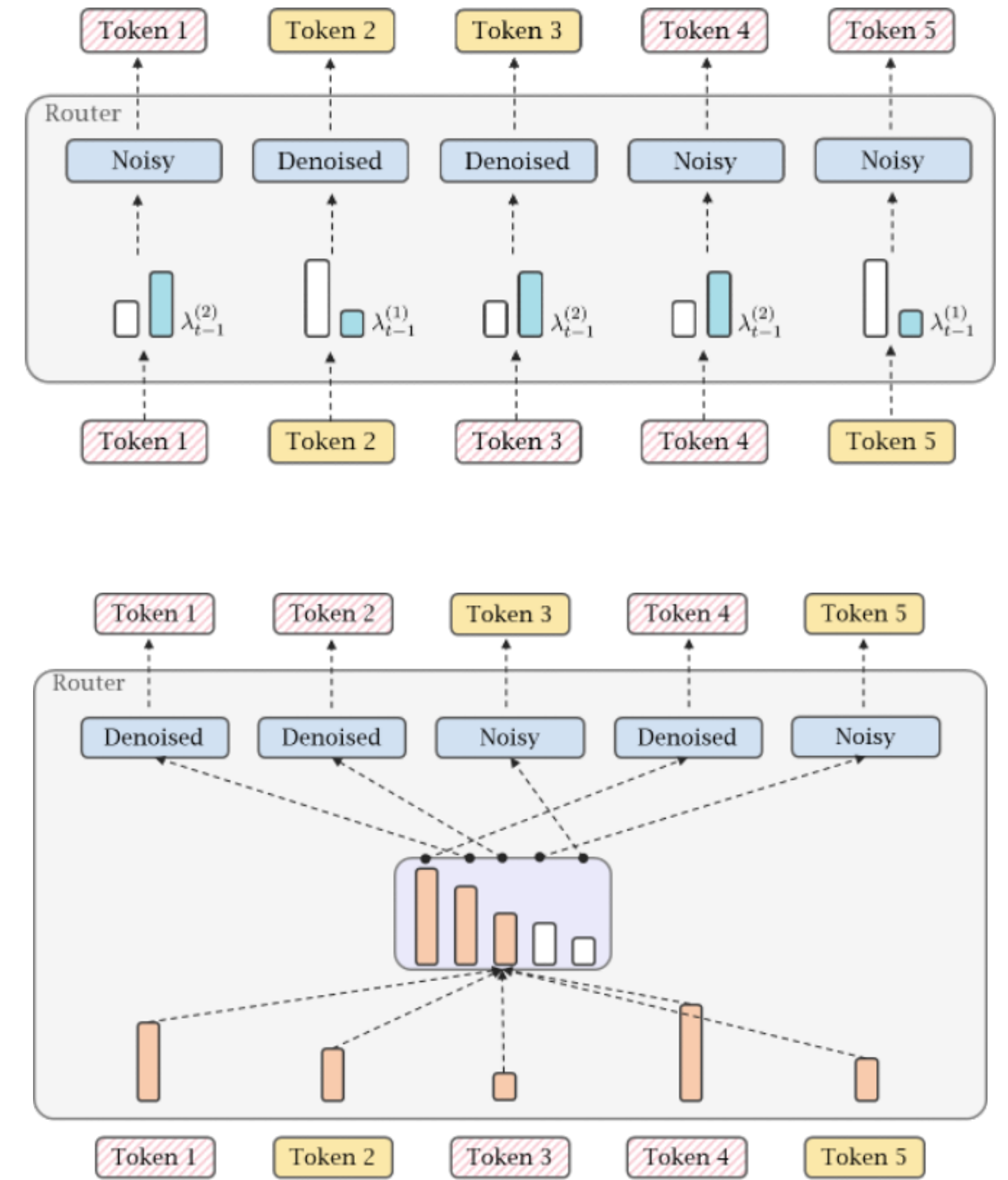
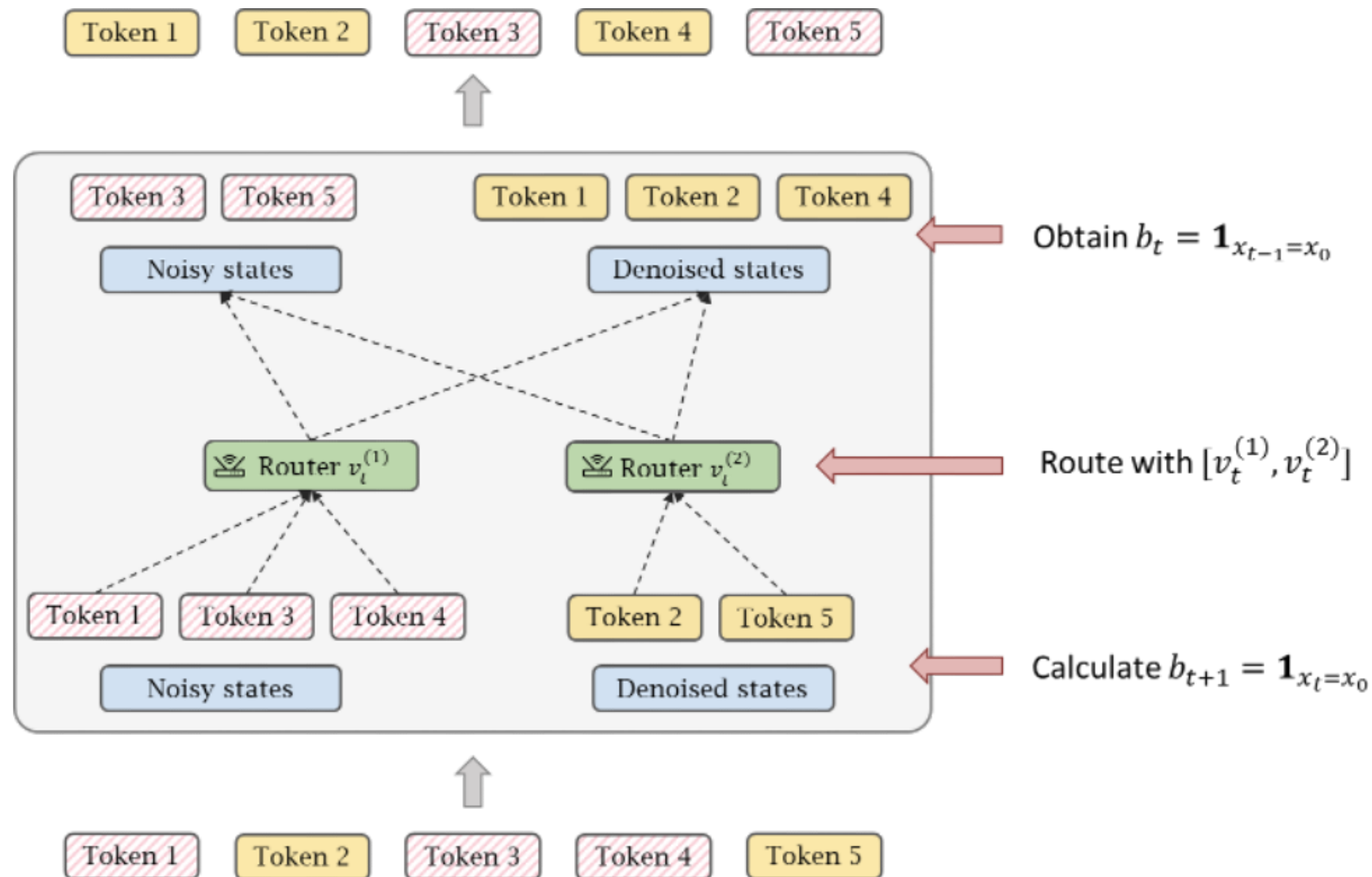
$$\lambda_{t-1}^{(2)} = 1 - \frac{t-1}{T}$$

Router-agnostic

Evaluated at noisy tokens

Simple Cross-Entropy Loss

Decoding RDMs



Confident-based Router

Experimental Results of RDMs

Close the gap with autoregressive models (Machine Translation)

	Model	# Iterations	IWSLT14 DE-EN		WMT16 EN-RO		WMT14 EN-DE	
			Vanilla	Reparam.	Vanilla	Reparam.	Vanilla	Reparam.
Continuous Diffusion	CDCD (Dieleman et al., 2022)	200	–		–		20.0*	
Discrete Diffusion	Multinomial Diffusion (Hoogeboom et al., 2021)	2	23.05	28.01	26.61	30.16	4.28	21.43
		4	24.24	30.57	27.81	31.70	4.31	24.05
		10	21.28	32.23	25.25	33.00	6.94	25.63
		16	20.59	32.58	24.36	33.11	6.07	25.64
		25	20.06	32.84	23.94	33.31	3.69	26.04
	Absorbing Diffusion (Austin et al., 2021)	2	25.24	27.60	27.24	30.72	16.46	21.00
		4	26.93	31.47	29.16	32.60	19.48	24.26
		10	28.32	33.91	30.41	33.38	21.62	26.96
		16	28.38	34.41	30.79	33.82	22.07	27.58
		25	28.93	34.49	30.56	33.99	22.52	27.59
Auto-regressive Models	Transformer-base (Vaswani et al., 2017)	n.a.	34.51		34.16		27.53	

Experimental Results of RDMs

Fix previously made errors in where vanilla discrete diffusion models get stuck

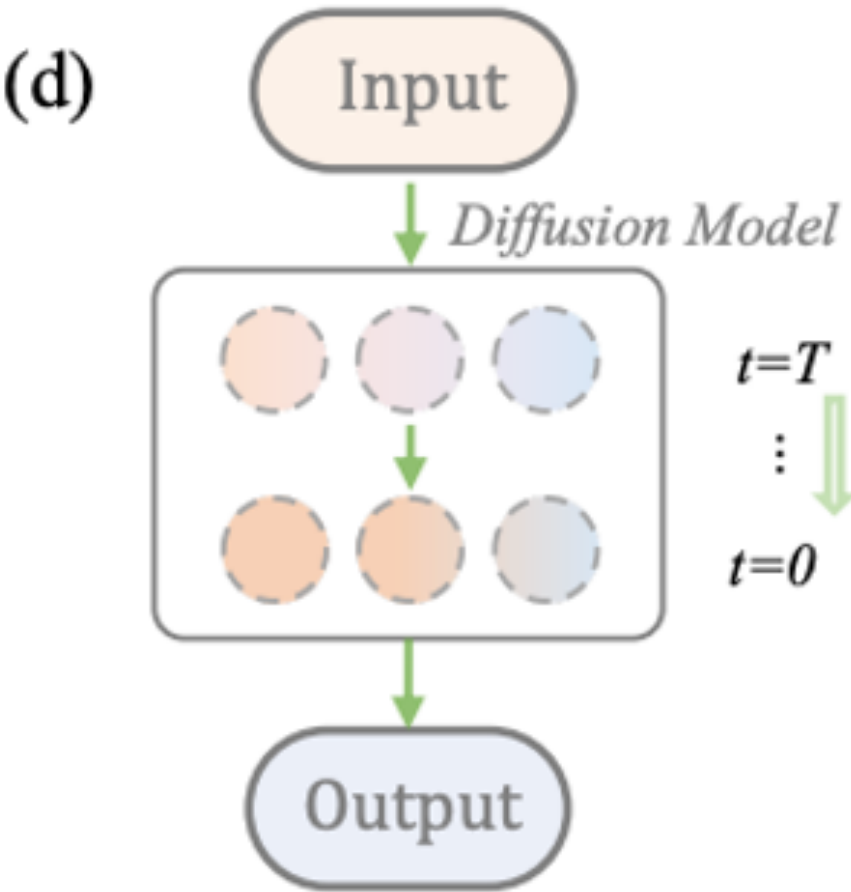
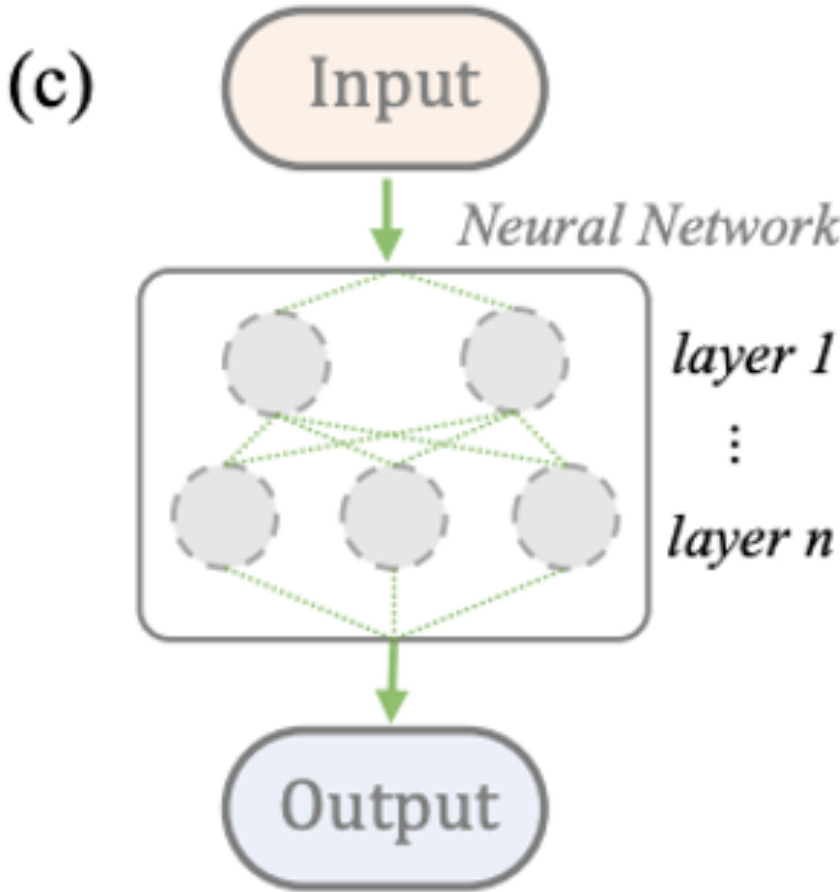
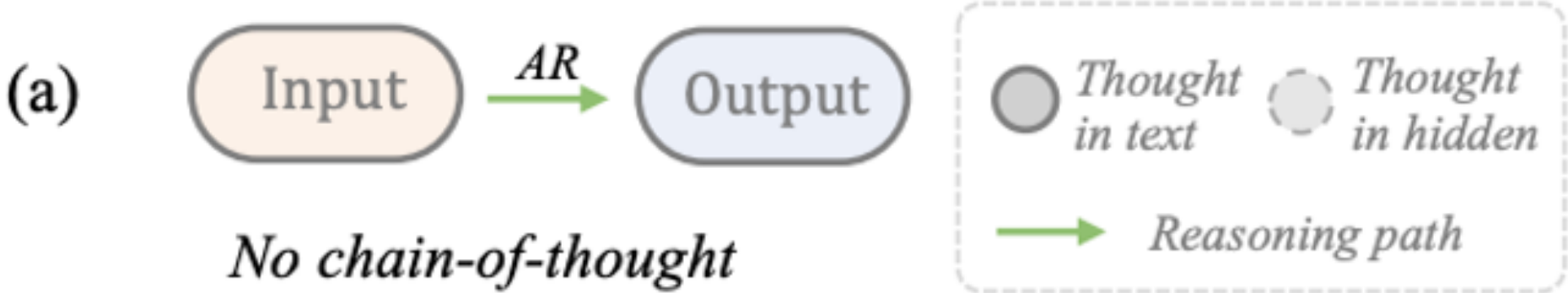
		Source: i have only 2 months for my ca cpt exams how do i prepare?
		Reference: i want to crack ca cpt in 2 months. how should i study?
		# Iter. Decodes
Absorbing	0	o <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M>
	1	o <M> <M> <M> <M> <M> ca <M> ca <M> <M> ca <M> <M> <M> <M>
	2	o <M> <M> <M> <M> <M> ca <M> ca <M> <M> ca <M> <M> <M> <M>
	3	o <M> <M> <M> <M> <M> ca - ca cp <M> ca <M> <M> exam <M>
	4	o <M> can <M> prepare <M> ca - ca cp <M> ca <M> ##t exam ?
	5	o how can i prepare for ca - ca cp ##t ca cp ##t exam ?
RDM-absorbing	0	o <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M>
	1	o how <M> i <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> <M> ?
	2	o how <M> i prepare for ca <M> ##t <M> <M> <M> <M> <M> ?
	3	o how <M> i prepare for ca cp ##t <M> <M> <M> months months ?
	4	o how <M> i prepare for ca cp ##t exam in two months <M> ?
	5	o how can i prepare for ca cp ##t exam in two months left ?

Autoregressive LLMs, what else?

Diffusion of Thoughts (DoTs)



Diffusion of Thoughts



Questions:

When Freda cooks canned tomatoes into sauce, they lose half their volume. Each 16 ounce can of tomatoes that she uses contains three tomatoes. Freda's last batch of tomato sauce made 32 ounces of sauce. How many tomatoes did Freda use?

CoT Decoding:

<

DoT Decoding:

<<32/16=2>> <<12/12=2*3 ## *2=2= #####

MP-DoT Decoding:

<<32==2>>

Multimodal LLMs

- Motivation: Why do we need a versatile large multimodal models (LMMs)?
- Progress: What recipe do we have for building capable LMMs?
 - Architecture
 - Data
 - Training Recipe

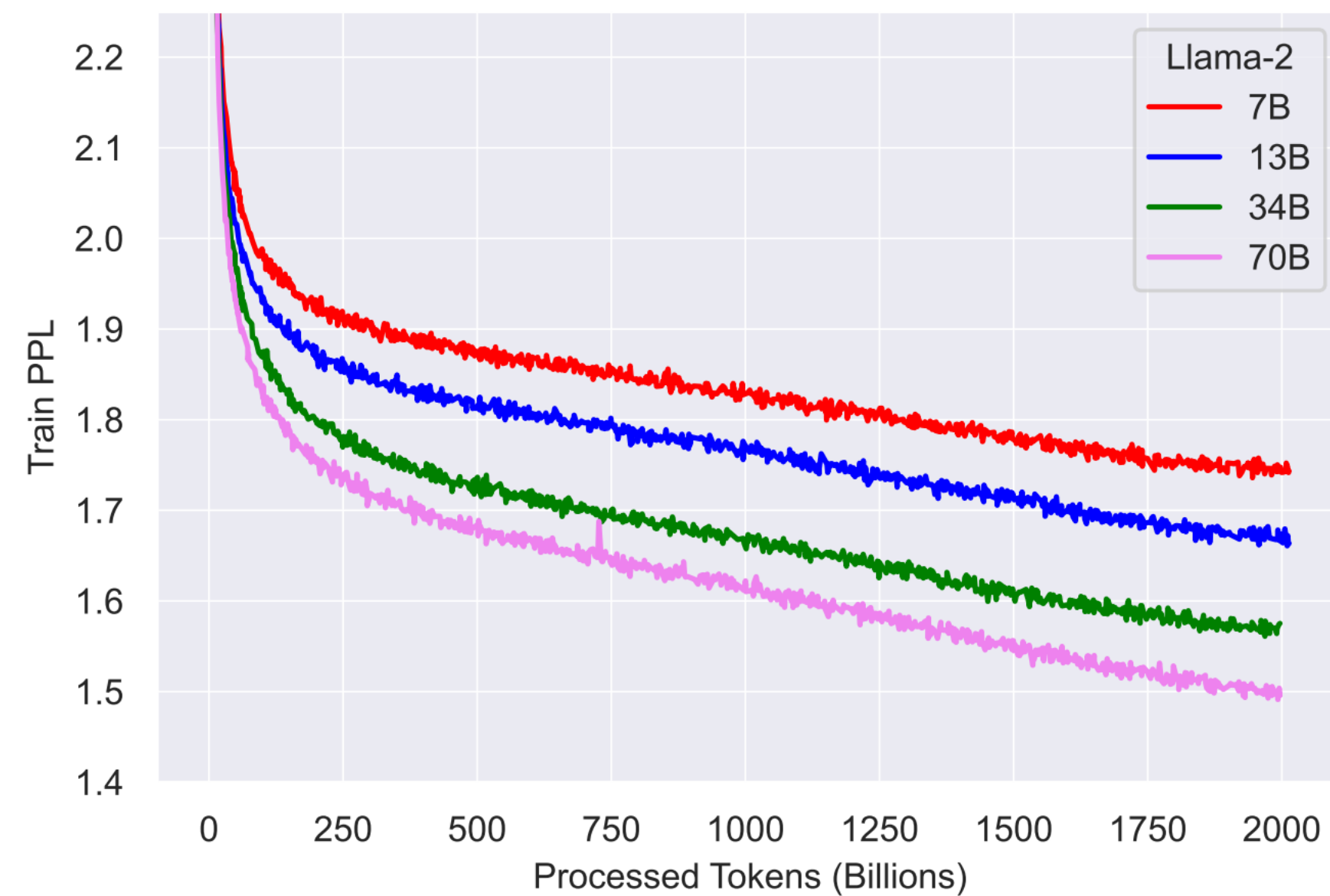
Future Directions: What we can do in the area of LMMs?



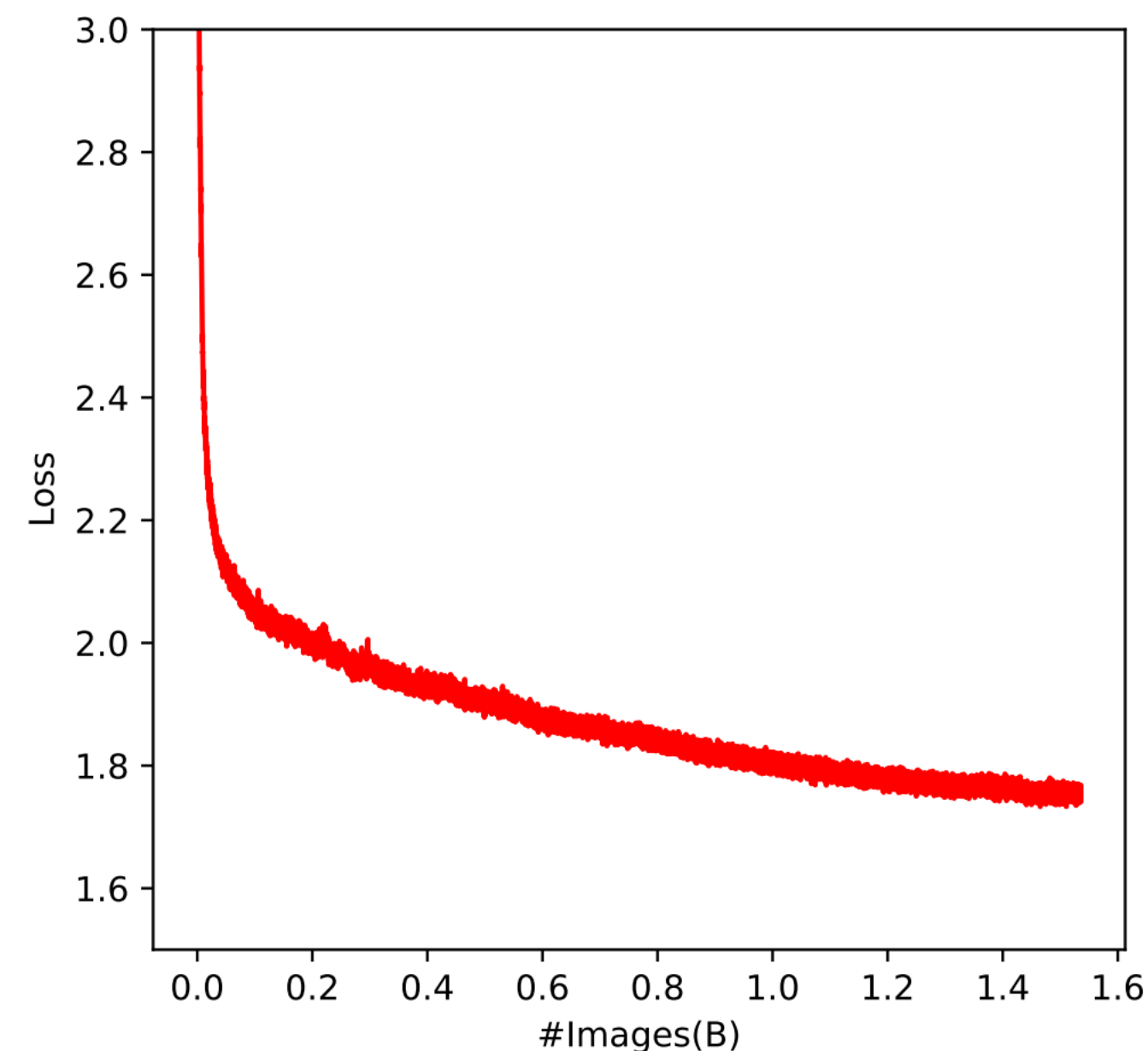
Slides for Multimodal LLMs from Lei Li (<https://lilei-nlp.github.io/>), xhs ID: tobiaslee

Motivations: Why LLMs

Compression of AGI



LLaMa-2 Training Loss



Qwen-VL-Chat (7B)

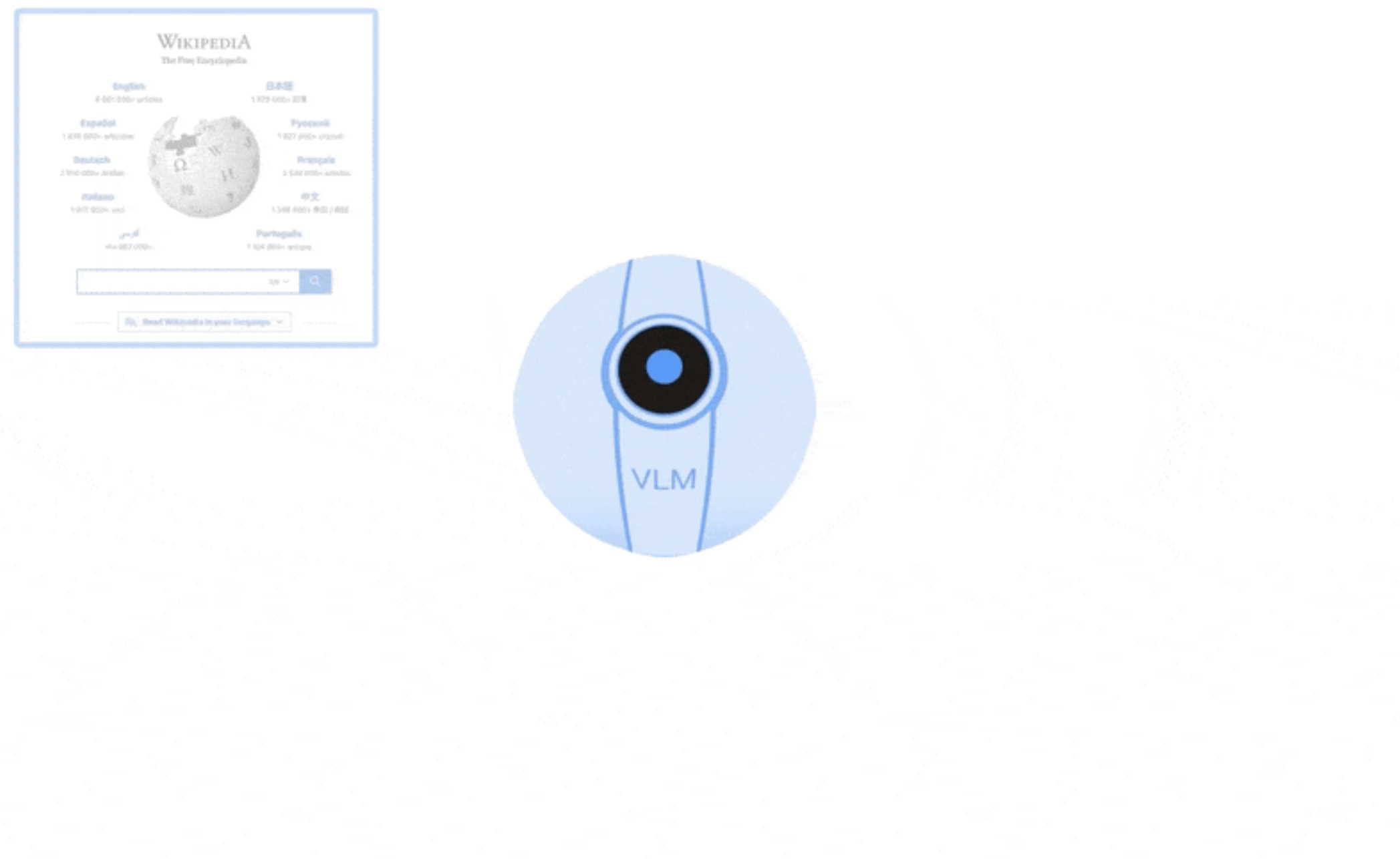
Image = 256 tokens
Caption \sim 64 tokens

\sim 500B tokens \rightarrow \sim 1.8

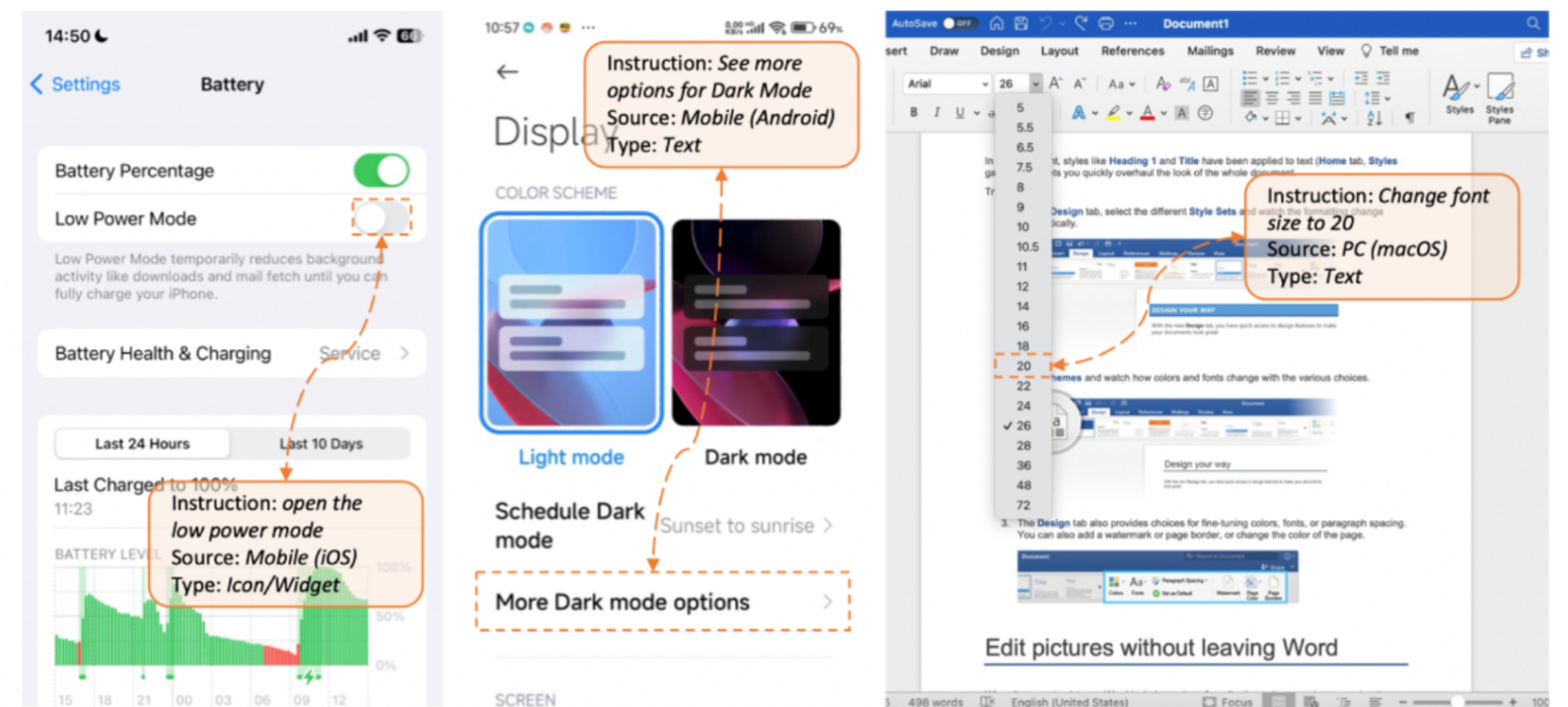
LLaMA-2-7B 500B tokens
 \rightarrow 1.9

Motivations: Why LLMs

Real-world applications are beyond text.



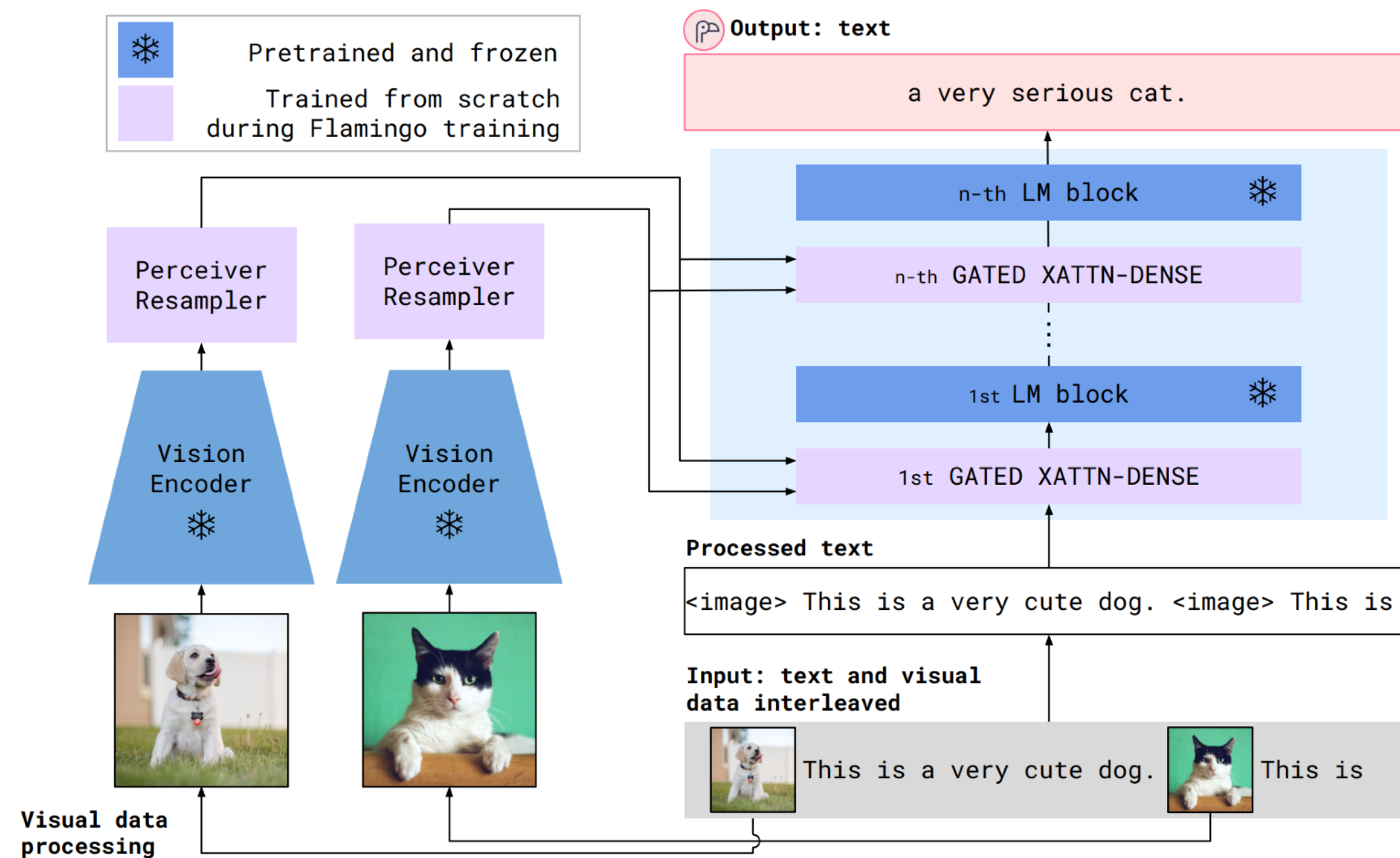
Embodied Robot



GUI Agent

RT-2: New model translates vision and language into action
SeeClick: Harnessing GUI Grounding for Advanced Visual GUI Agents

Architecture Overview



Flamingo

Images are integrated via
Gated Cross-Attention

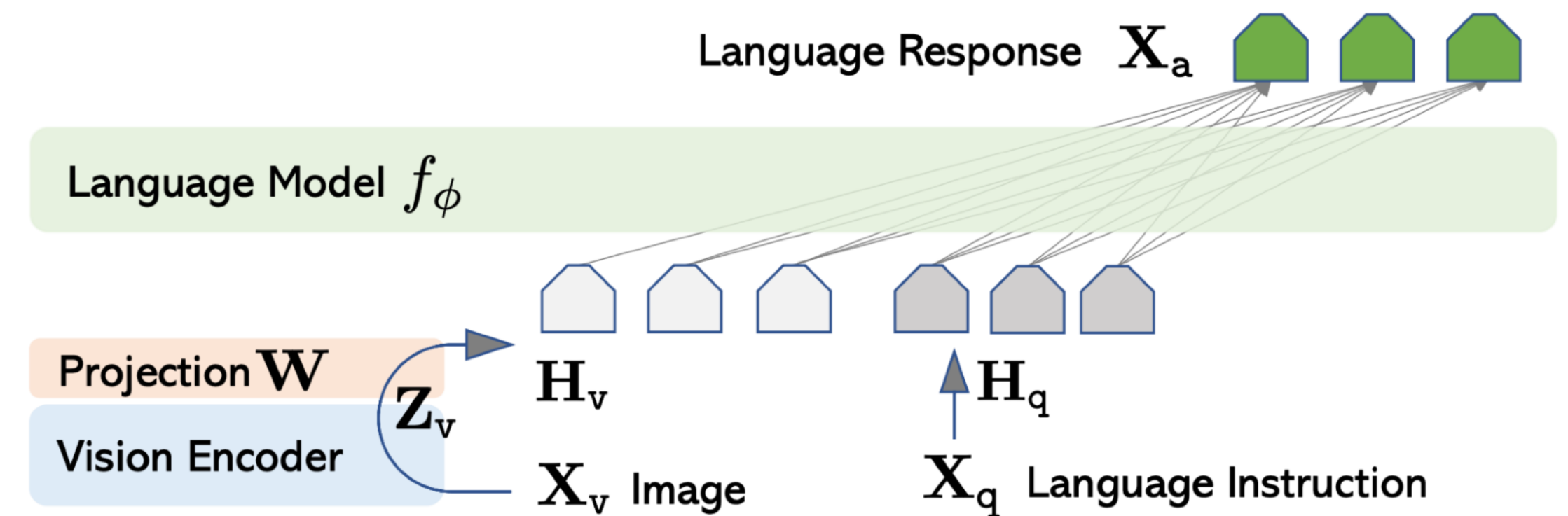


Figure 1: LLaVA network architecture.

LLaVA / Qwen-VL-Chat / MiniGPT4 ..

Images are treated as Word Embeddings
More prevailing now

Architecture Design Space

- Large Language Model:
 - Use the best LLM available.
- Vision Encoder:
 - Influence of the image resolution.
 - Which vision encoder shall we use?
- Modality Projector:
 - How to **effectively** represent the image in word embedding space?

How might GPT-4v deal with high resolution?

Here are some examples demonstrating the above.

- A 1024×1024 square image in `detail: high` mode costs 765 tokens
 - 1024 is less than 2048, so there is no initial resize.
 - The shortest side is 1024, so we scale the image down to 768×768 .
 - 4 512px square tiles are needed to represent the image, so the final token cost is $170 * 4 + 85 = 765$.
- A 2048×4096 image in `detail: high` mode costs 1105 tokens
 - We scale down the image to 1024×2048 to fit within the 2048 square.
 - The shortest side is 1024, so we further scale down to 768×1536 .
 - 6 512px tiles are needed, so the final token cost is $170 * 6 + 85 = 1105$.
- A 4096×8192 image in `detail: low` mode costs 85 tokens
 - Regardless of input size, low detail images are a fixed cost.

GPT-4V Document

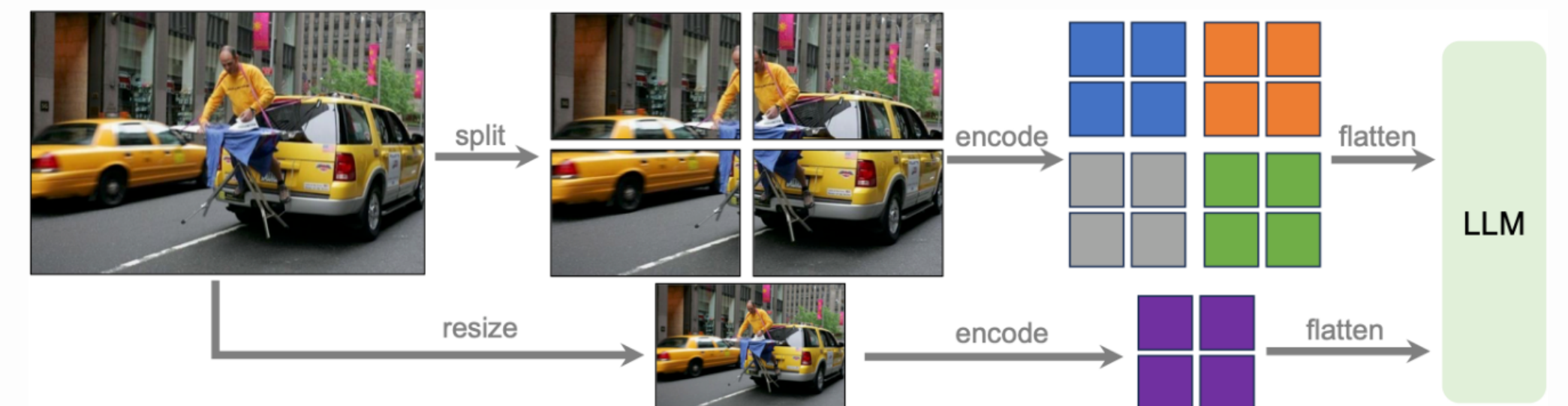
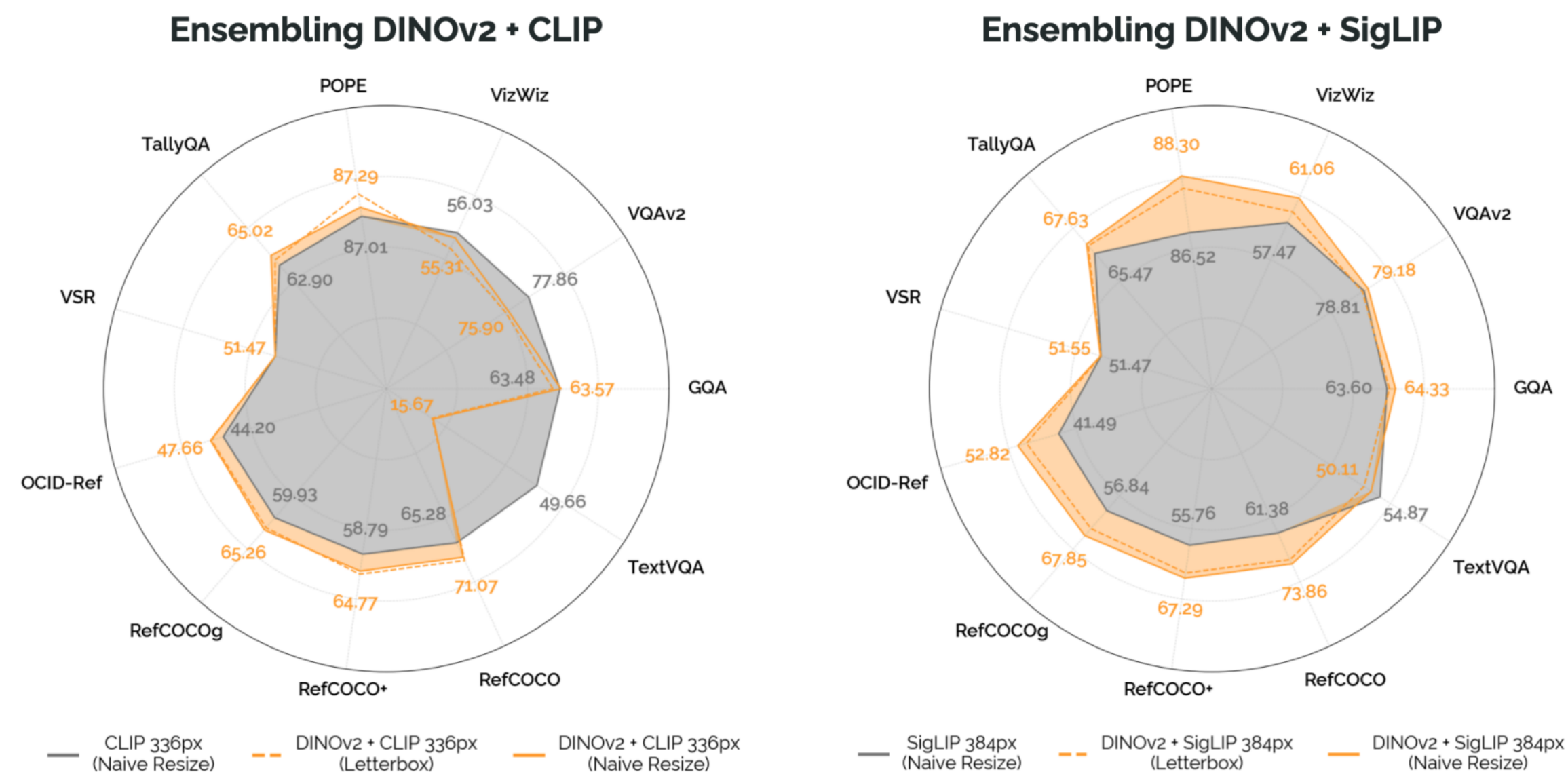


Illustration of dynamic high resolution scheme: a grid configuration of 2×2

LLaVA-Next cuts image into grids to reduce the computation overhead

LLaVA-NeXT: Improved reasoning, OCR, and world knowledge

Vision Feature Types



- Single CLIP Feature is not comprehensive for grounding related tasks.
- Combination of SigCLIP and DINOv2 (image SSL) leads to the best results across benchmarks.

GPT-4V Document

Modality Projector

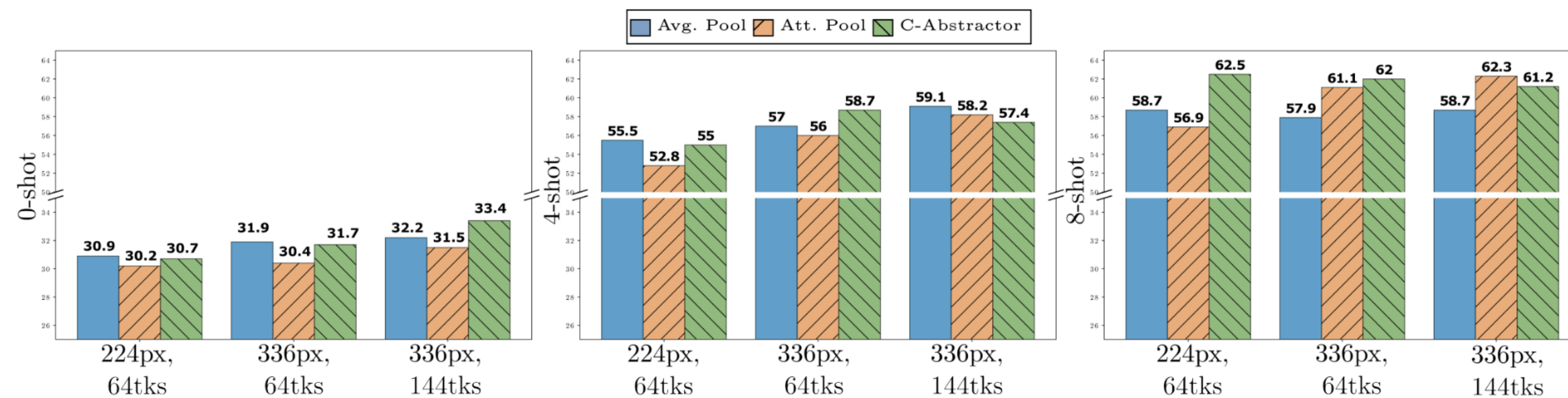
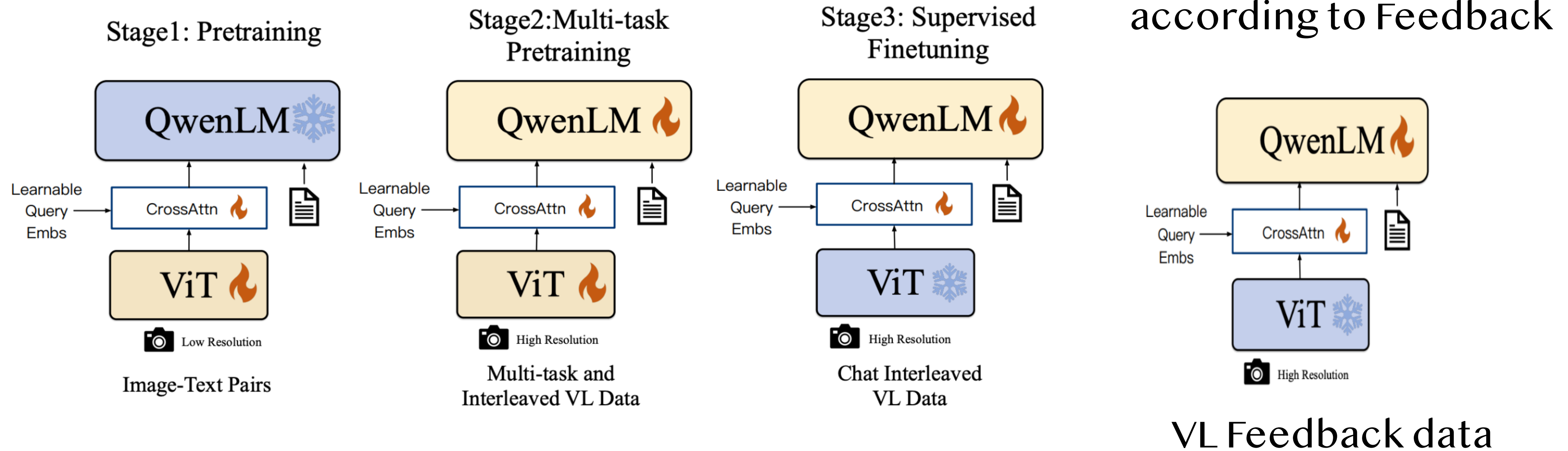


Fig. 4: 0-shot, 4-shot, and 8-shot ablations across different visual-language connectors for two image resolutions, and two image token sizes.

- Projector can be utilized to compress the visual tokens as well, which one shall we use?
- Avg Pooling v.s. Attn Pooling v.s C-Abstractor (ResNet Pooling)
- Findings:
 - Number of visual tokens and image resolution matters most, while the type of VL connector has little effect.

Data & Training Stage Overview



General VL Datasets

- Stage 1: Large-scale Paired dataset for modality alignment
 - Scale & Quality
 - Raw datasets: Laion5B, Datacomp-1B
 - LLM (LMMs) rewrites dataset: Laion-COCO, Capsfusion, ShareGPT4V
- Stage 2 & 3: Instruction tuning datasets / Chatty Dataset
 - M3IT (Academic VL tasks, <https://m3-it.github.io>), wide coverage.
 - LLaVA-Instruct (GPT-4/ChatGPT generated pseudo-multimodal dataset)
 - Specialized Domain Datasets: AI2D, OCR-VQA, LLaVAR, ChartVQA, G-LLaVA



COCO: Young children standing on a platform waiting for a train to arrive. Adults and children watching a train slowly leave. A family near a railroad track watching the train pass. People waiting on a platform as a train pulls up. A train station with a green train on the tracks and children waiting for it to go by.

LLaVA: At a train station, a group of people, including both young children and adults, are standing on a platform waiting for a train to arrive. The train is already present on the tracks, partially visible on the right side of the image. Some of the people watch the train closely, while others seem to be patiently anticipating its departure.

There is a total of eight individuals waiting for the train, with one child in the middle of the platform and the others scattered around. A backpack can be found on the far left side of the platform, suggesting that someone may have set it down while waiting.

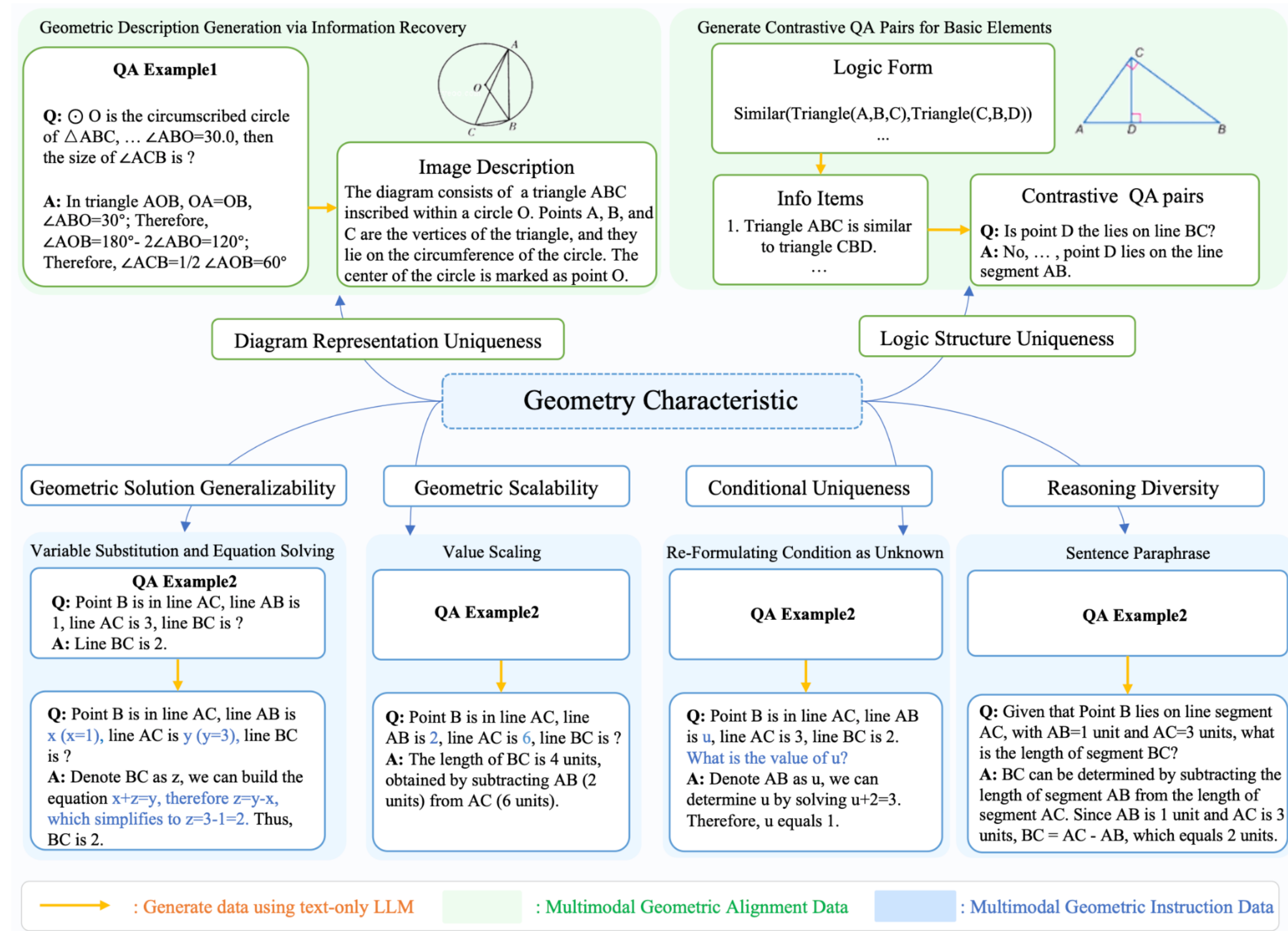
ShareGPT4V: The image captures a moment at a train station. A green train is on the tracks, moving away from the platform labeled "Platform 2". The train's motion is observed by people standing on the platform, their attention drawn to the departing vehicle.

A red signal stands on the tracks, its vibrant color contrasting with the green of the train. Signs are prominently displayed around the platform. One warns "Beware of Trains", a cautionary reminder of the powerful machines that frequent these tracks. Another sign advises "Please Use The Footbridge To Cross The Line. If You Require Assistance Please Ask The Station Staff", guiding passengers to safely navigate the station.

The sky above is a clear blue, providing a serene backdrop to the bustling activity of the station. In the distance, trees can be seen, adding a touch of nature to this man-made setting. The image is a snapshot of everyday life at a train station, capturing both its routine operations and its inherent charm.

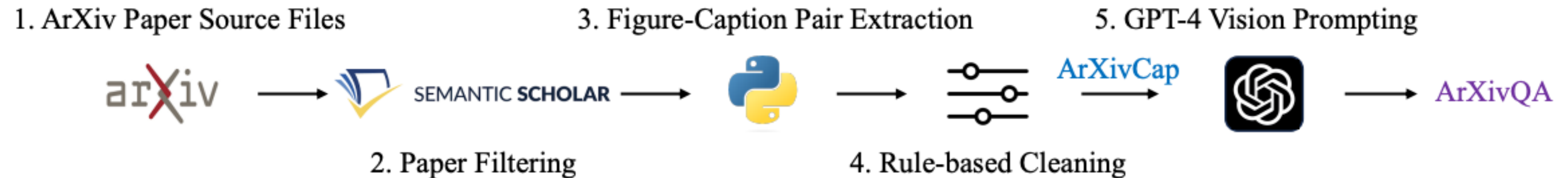
(a) Comparison of Captions' Quality

Datasets for LLMs Math Reasoning



G-LLaVA: Solving Geometric Problem with Multi-Modal Large Language Model

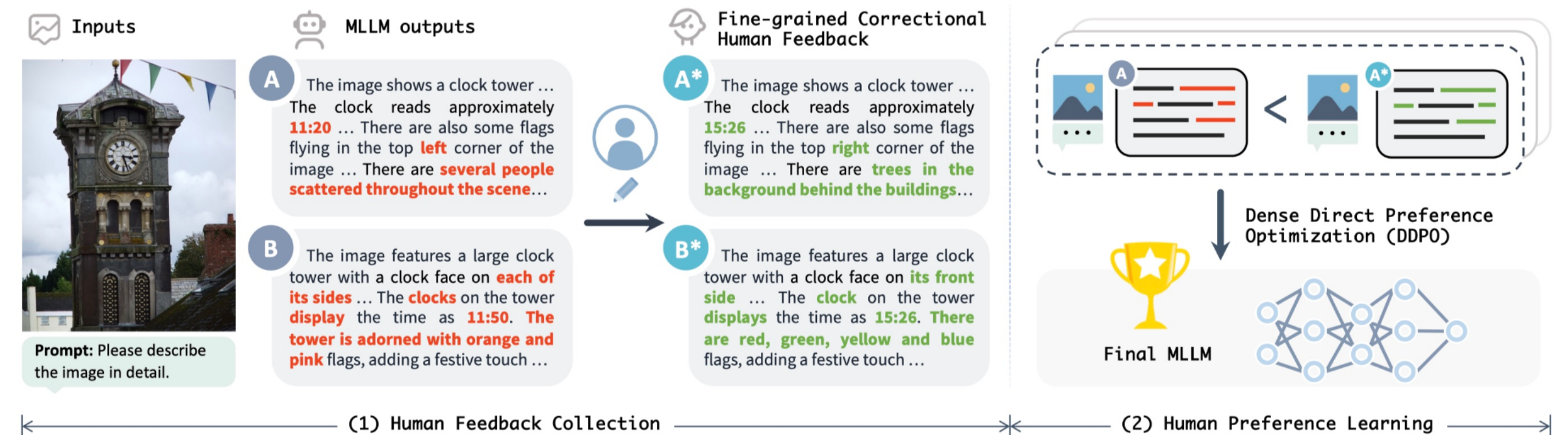
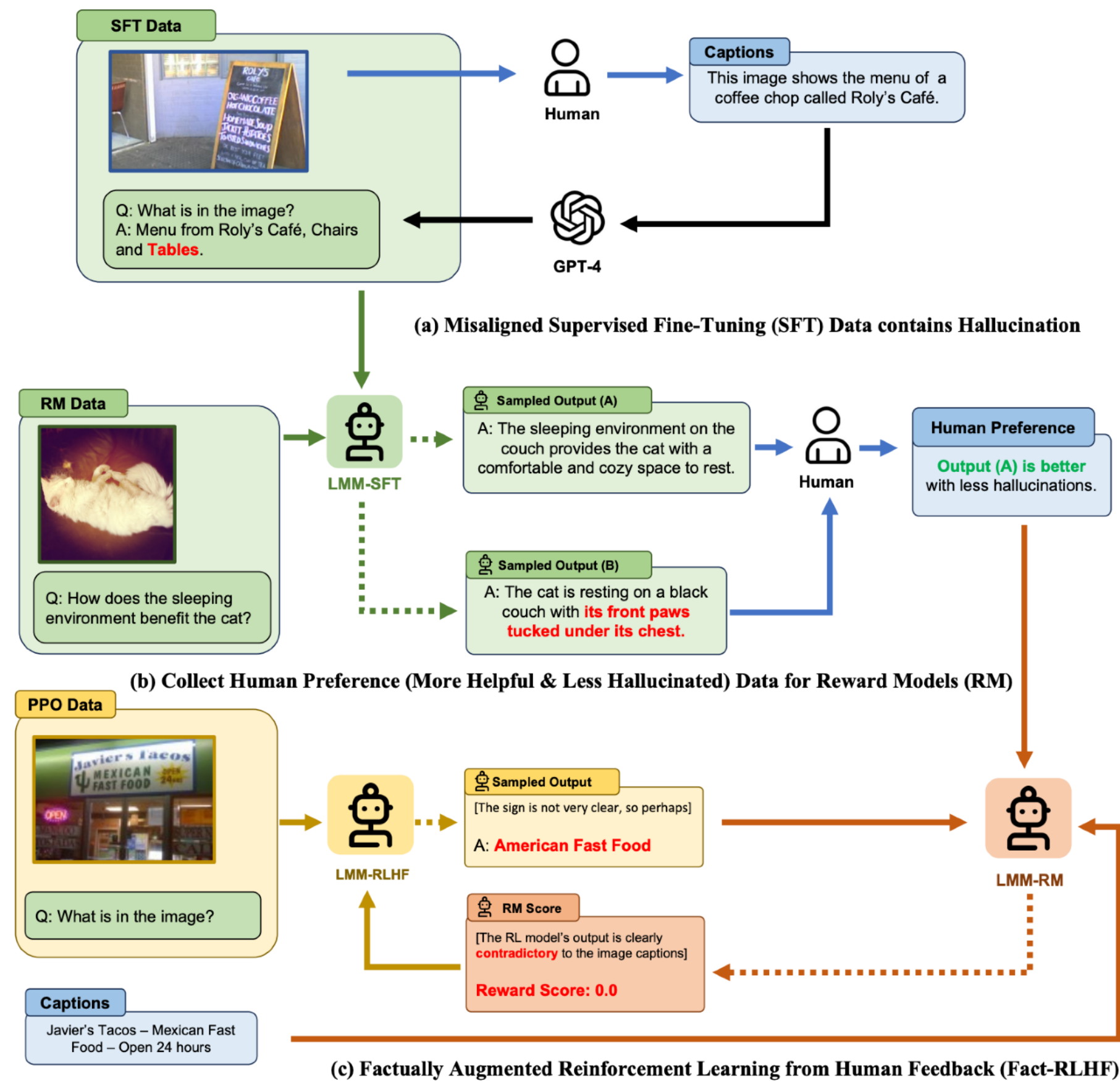
Multimodal ArXiv for LLMs Scientific Understanding



Dataset	Image Number	Paper Number	Image Category	Domain	Real Data
FigCAP (Chen et al., 2020)	219K	N / A	Bar, Line and Pie Charts	N / A	✗
SciCap (Yang et al., 2023b)	2.1M	295K	Open-Category	Computer Science and Machine Learning	✓
M-Paper (Hu et al., 2023)	350K	48K	Open-Category	Mainly "Deep Learning"	✓
ArXivCap (Ours)	6.4M	572K	Open-Category	Open-Domain	✓
FigureQA (Kahou et al., 2017)	140K	N / A	Bar, Line and Pie Charts	N / A	✗
DVQA (Kafle et al., 2018)	300K	N / A	Bar Charts	N / A	✗
ArXivQA (Ours)	32K	16.6K	Open-Category	Open-Domain	✓

Multimodal ArXiv: A Dataset for Improving Scientific Comprehension of Large Vision-Language Models

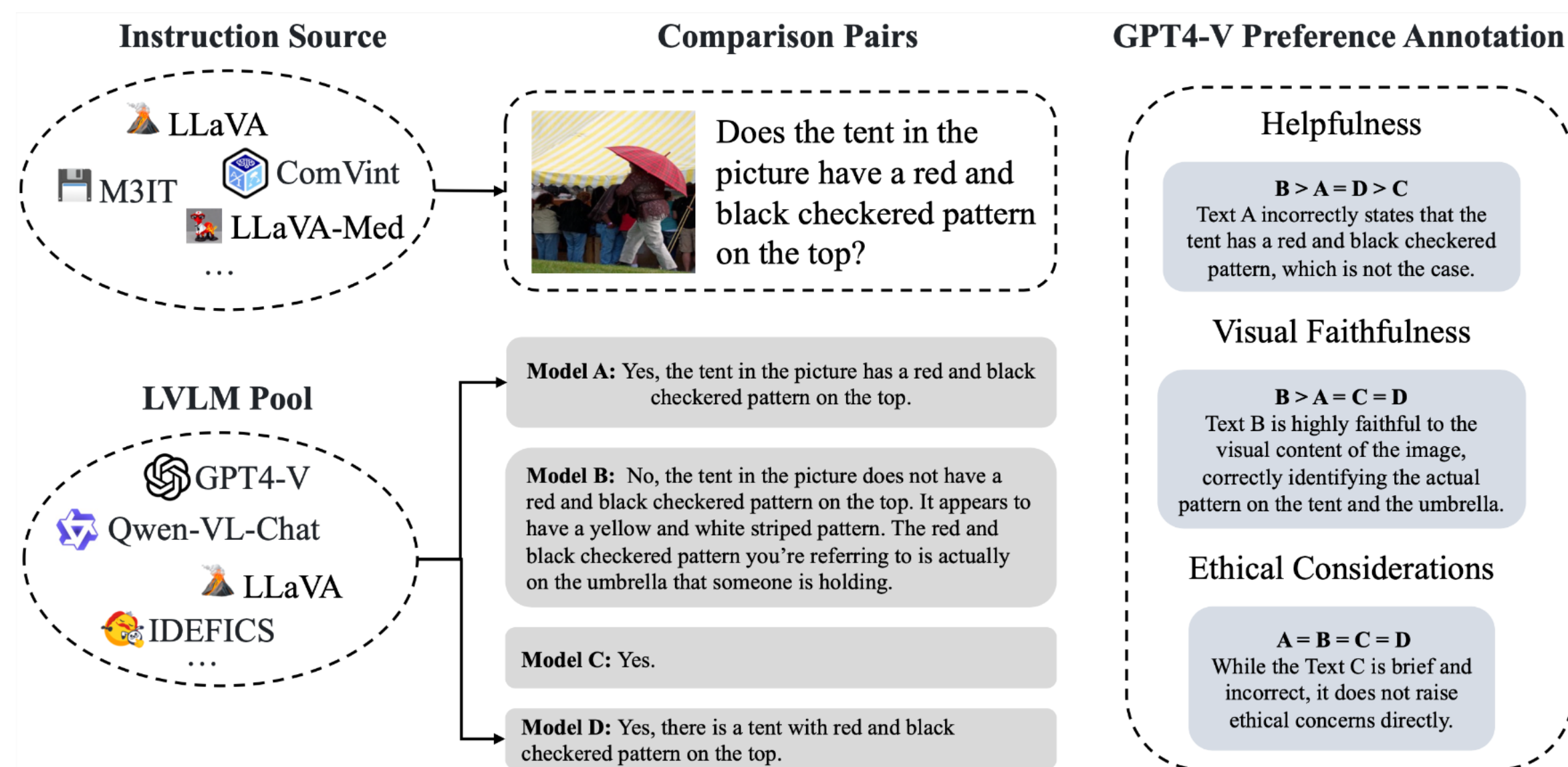
Feedback Dataset for LLMs



Aligning Large Multimodal Models with Factually Augmented RLHF

RLHF-V: Towards Trustworthy MLLMs via Behavior Alignment from Fine-grained Correctional Human Feedback

VLFeedback: Preference Distillation for LLMs



- We build a GPT-4V annotated vision-language feedback dataset, consists of 80K instructions decoded by various models, regarding three aspects.

A Guess of GPT-4V Recipe

- ChatGPT/GPT-4 as backbone LLMs for instruction following & reasoning
- Close-sourced Vision Encoders with 512px with tiling mechanism
 - Note CLIP is also from OpenAI
- Billion Scale (?) ChatGPT paraphrased detailed image-caption pairs for alignment training (inferred from DALL-E & SORA practice)
- Diverse vision-language tasks collected for SFT
- OCR pipeline integrated (?)
- Professional alignment team for eliminating bias