

Can Generative Al and Video Analysis Redefine the Future Vision? Unveiling the Potential and Challenges

Nick Tai, NVIDIA AI Technology Center (NVAITC)

2023/12/11



Introduction

Video Action Anticipation

Introduction to Diffusion Models

Synthesis and Future Directions

Conclusion and Discussion



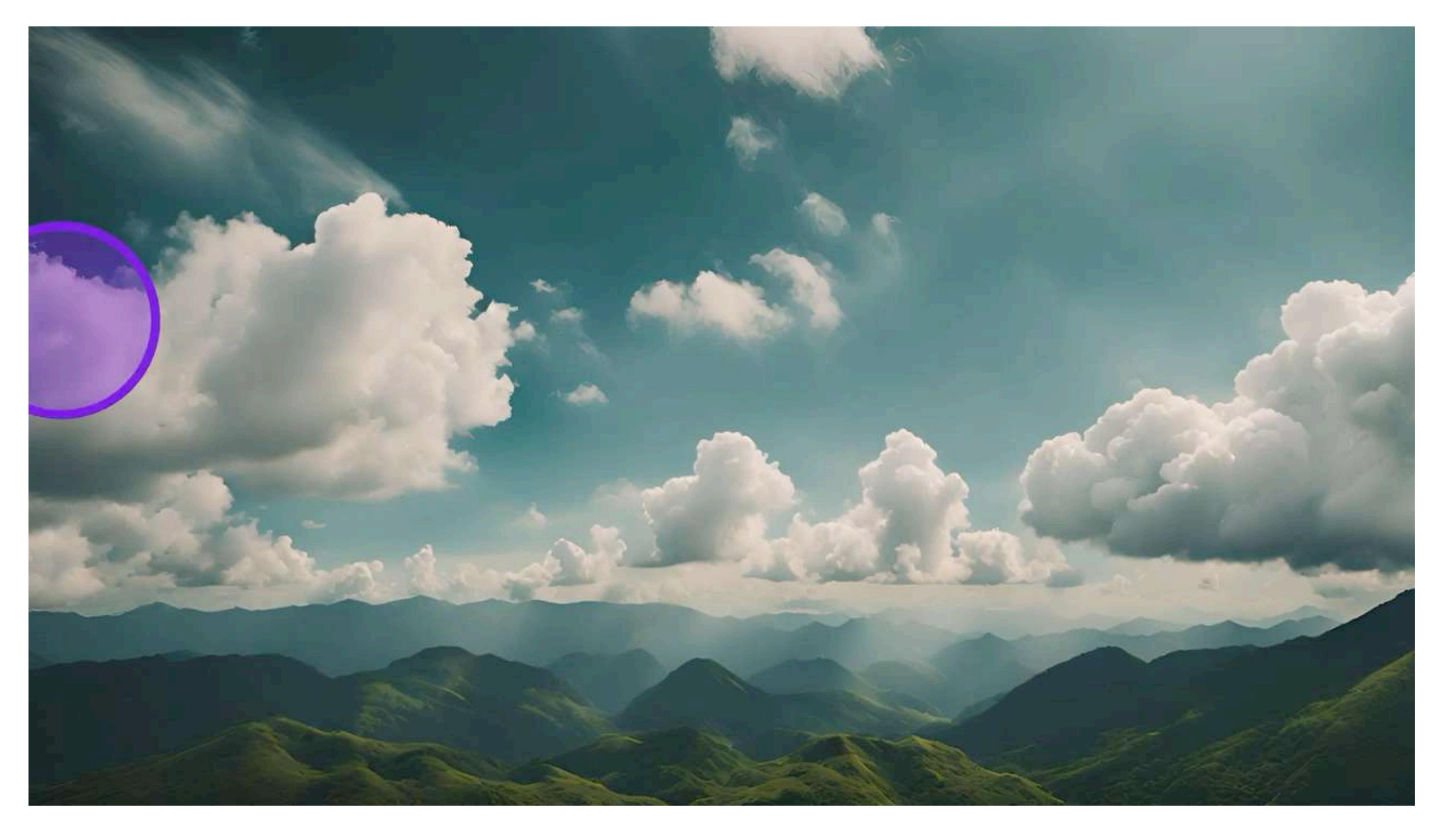
Introduction

Video Action Anticipation

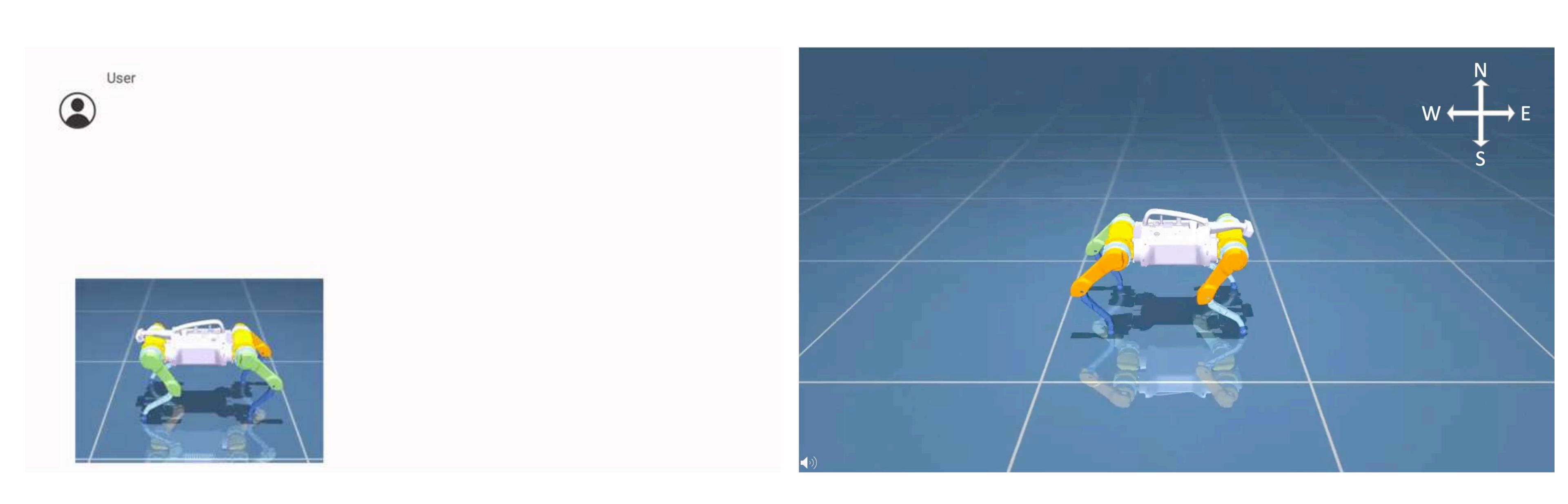
Introduction to Diffusion Models

Synthesis and Future Directions

Conclusion and Discussion

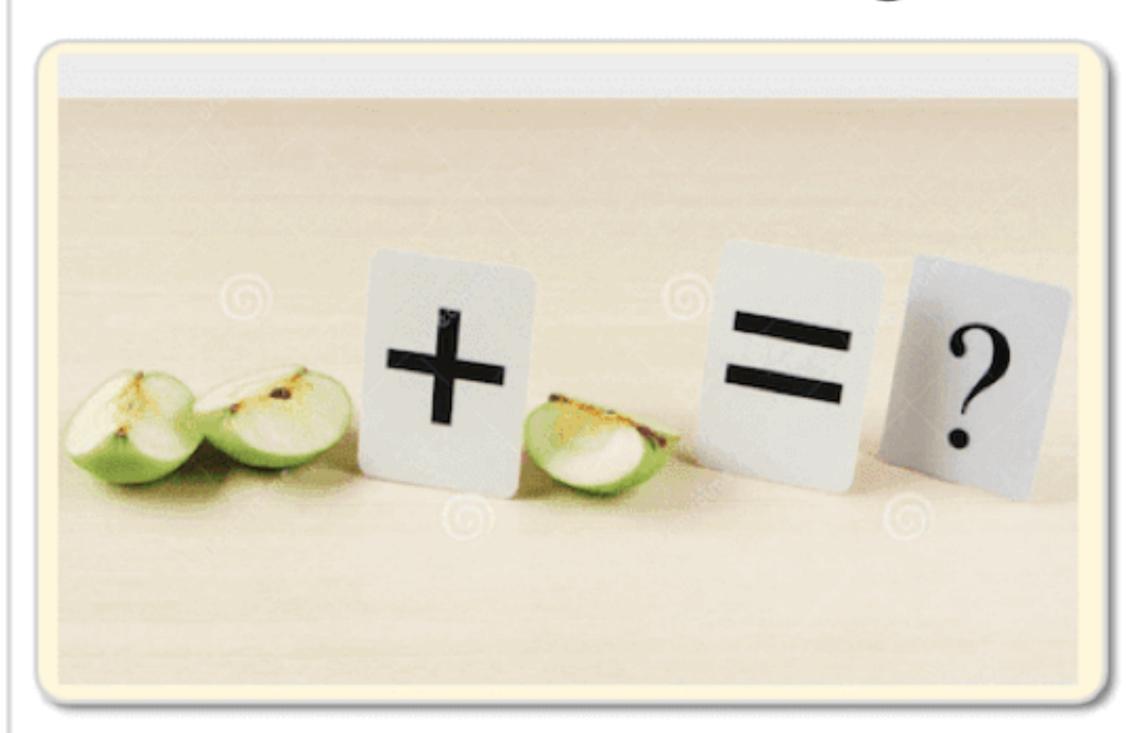


GEN-2 (Runway, 2023)



Language to Rewards for Robotic Skill Synthesis (Yu et al., 2023)

Visual Math Problem Solving



Q1: How to solve it?

CogVLM: The problem in the picture is a simple math problem, where a sliced green apple is equated to a question mark. Specifically, in the middle of the picture, there are three green apples placed side by side and three cardboard signs. The signs have black characters on them: a plus sign, an equal sign, and a question mark. The equation shown in the picture is: 2 + 1 = ?

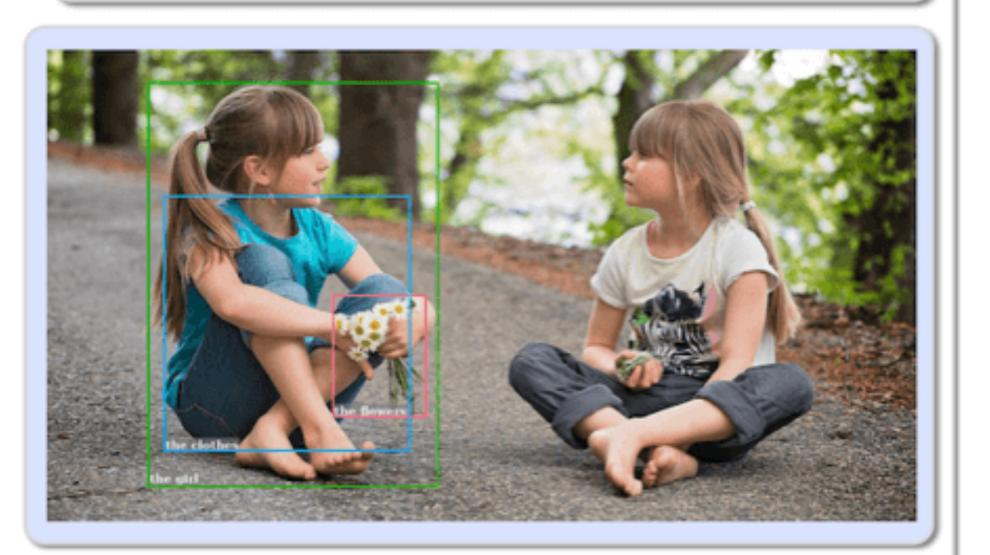
Q2: What's the final answer?

CogVLM: 2+1=3

Grounding Visual Question Answering

Q1: What color are the clothes of the girl whose hands are holding flowers? Let's think step by step.

CogVLM: Find the flowers [[328,519,436,780]]. Check the girl holding it, got the girl [[116,069,450,927]]. Check the clothes of it, got the clothes [[134,309,418,852]]. The clothes [[134,309,418,852]] are blue. So the answer is blue.

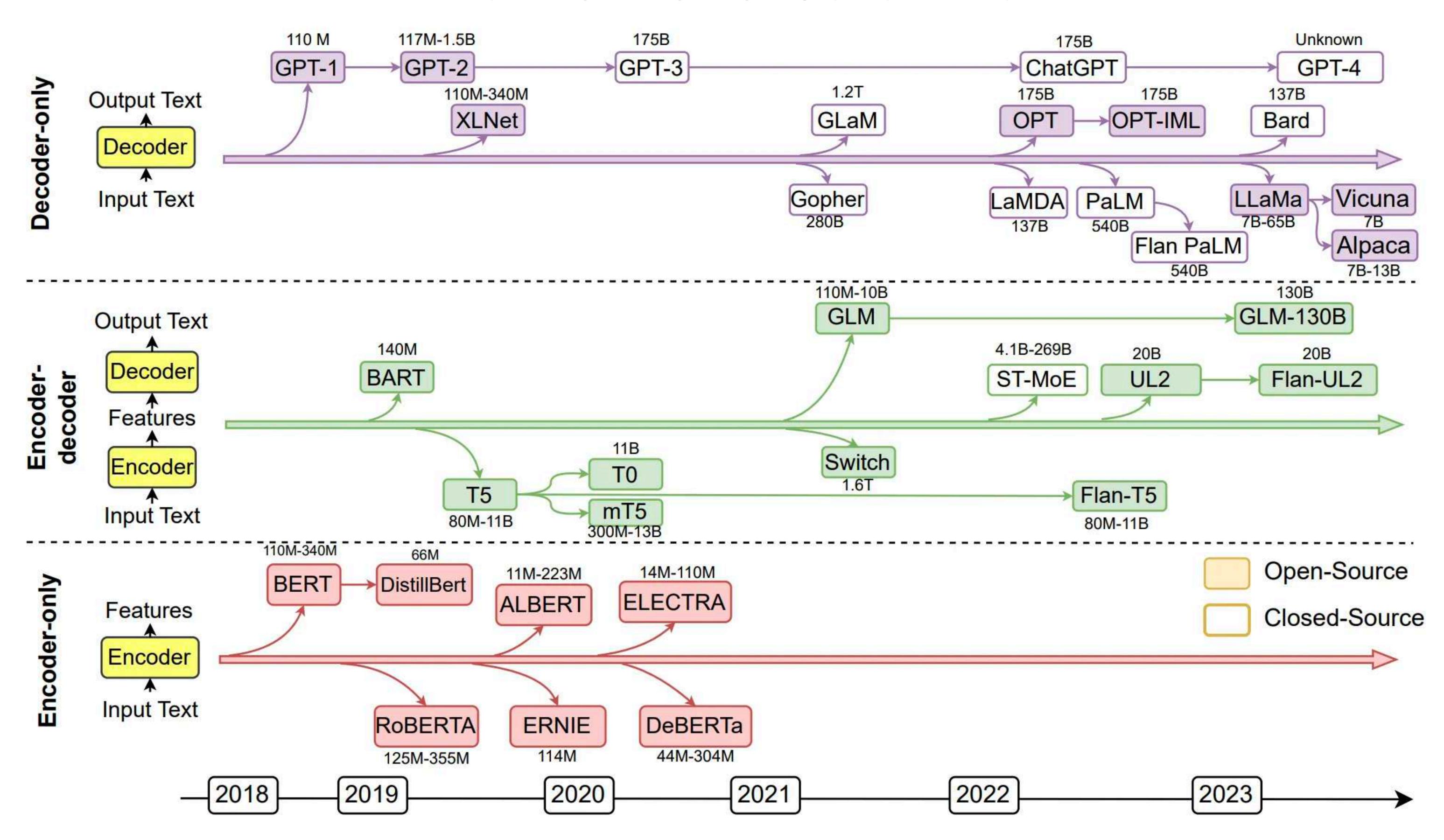


Q2: What color are the jeans of the girl whose hands are not holding flowers? Let's think step by step.

CogVLM: Find the not(flowers) hands [[654,637,710,728]]. Check the girl with it, got the girl [[526,109,880,946]]. Check the jeans of it, got the jeans [[530,613,880,871]]. The jeans [[530,613,880,871]] are gray. So the answer is gray.

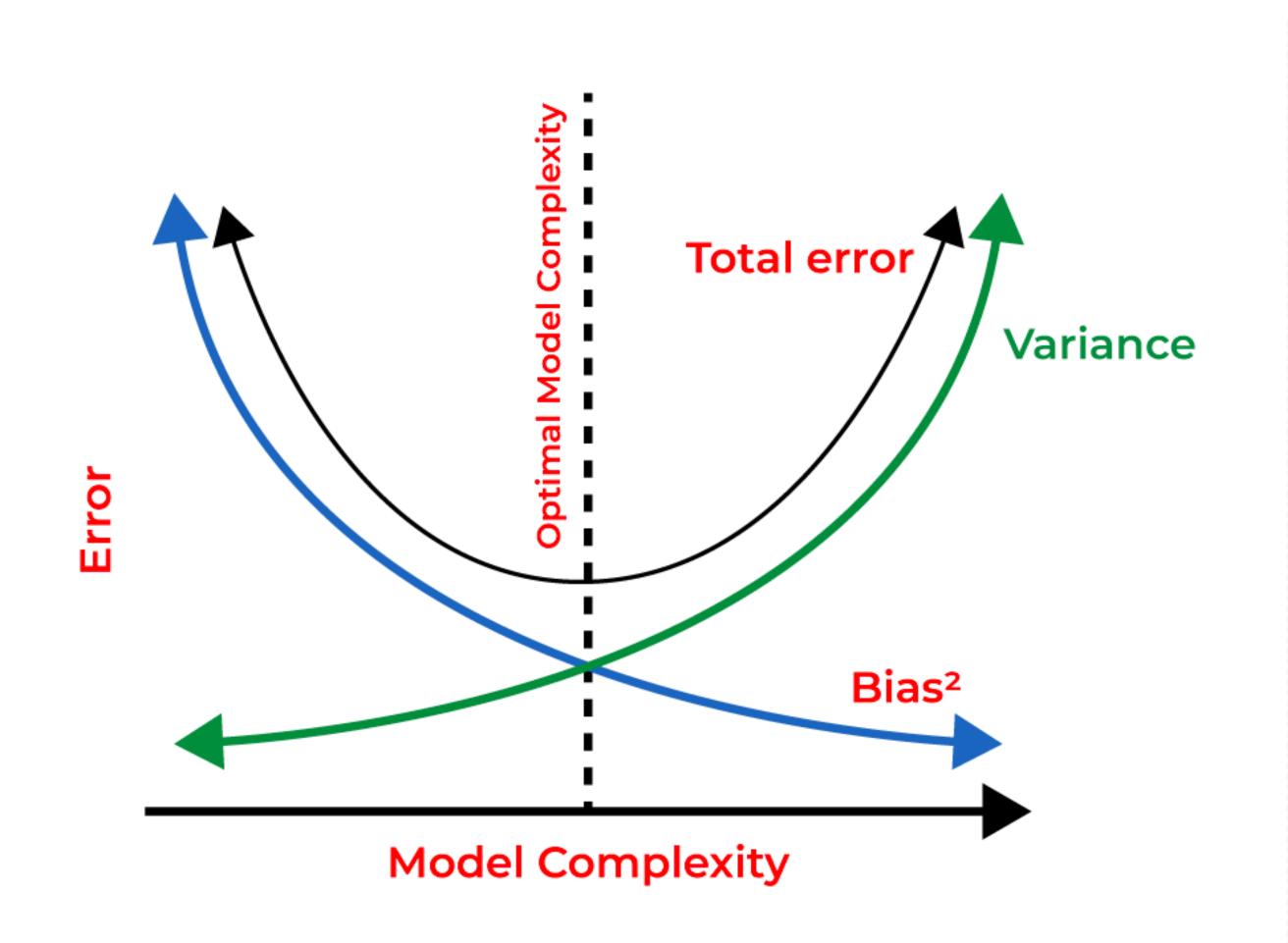


The Evolution of Generative Al

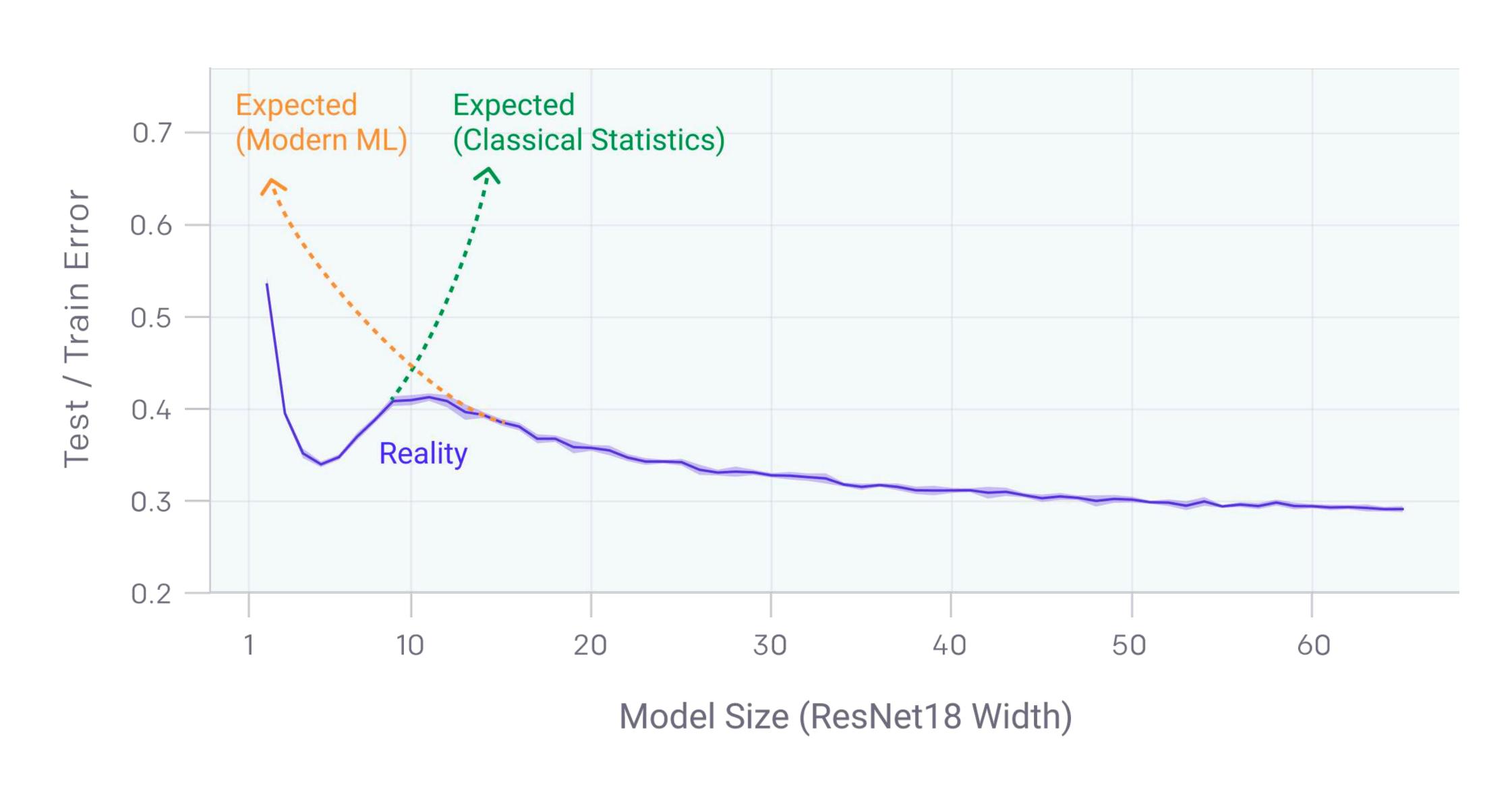


Why Large Models?

Double Decent Theory; (OpenAl 2019: https://openai.com/research/deep-double-descent)



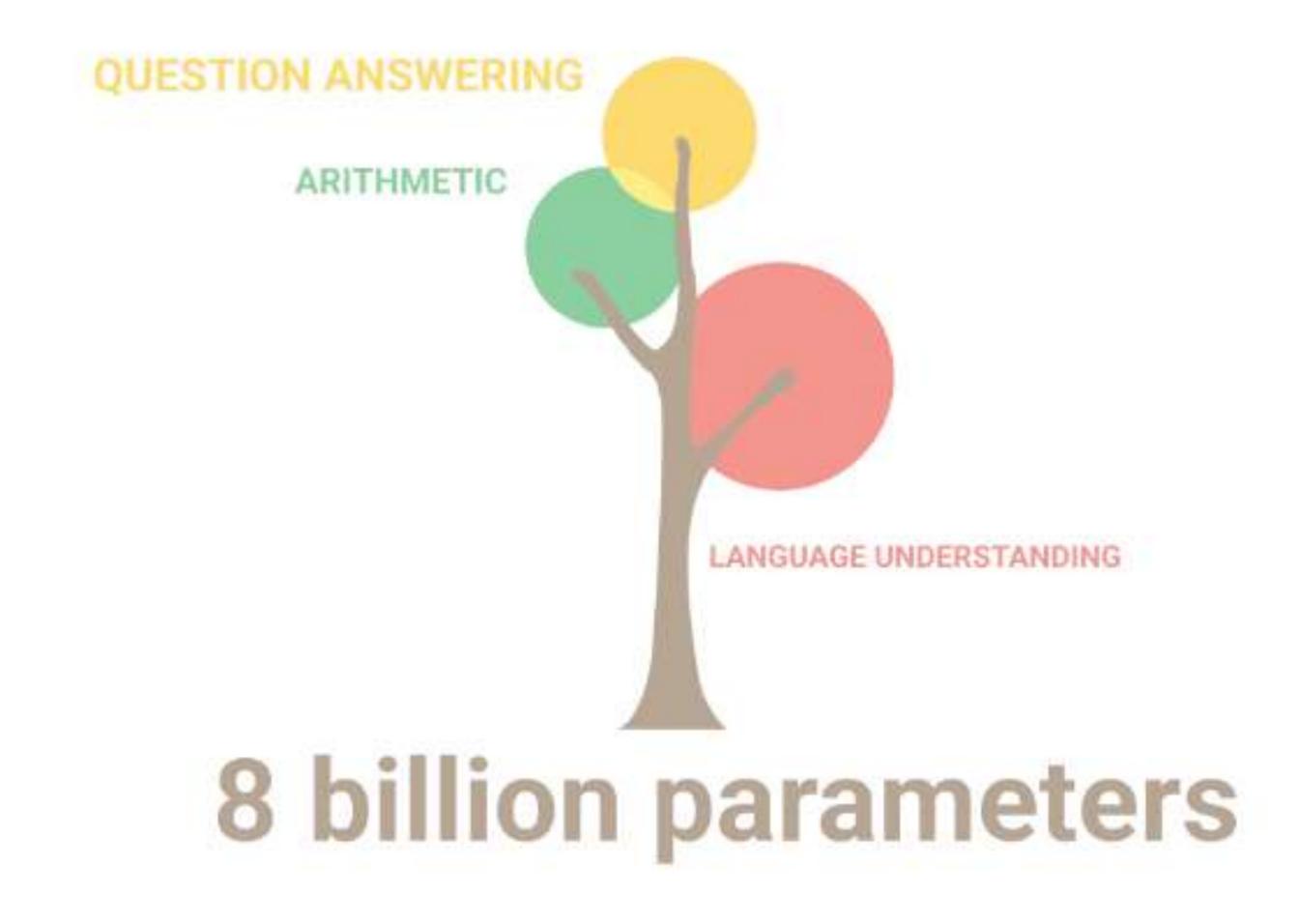






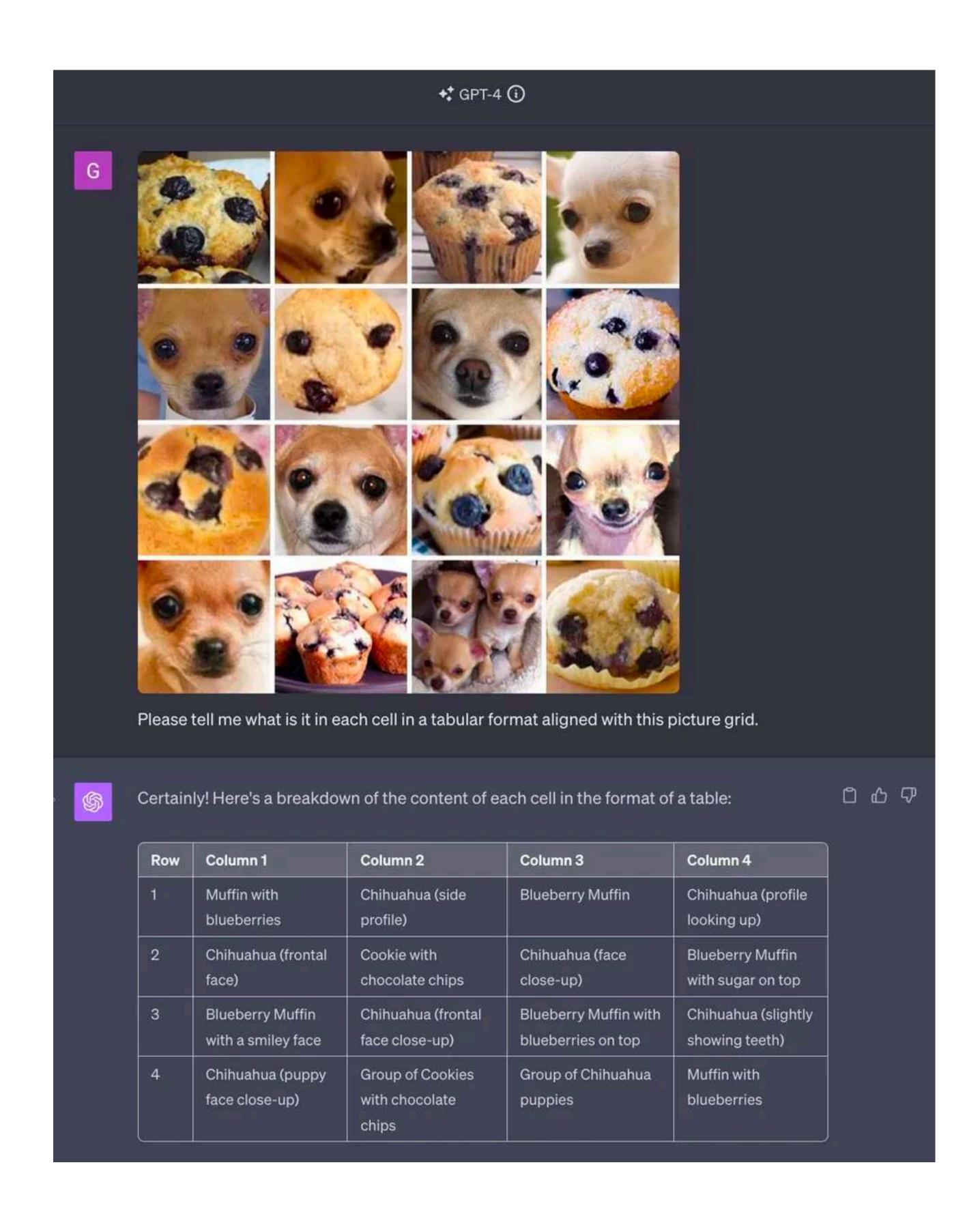
Why Large Models?

More weights, More strength



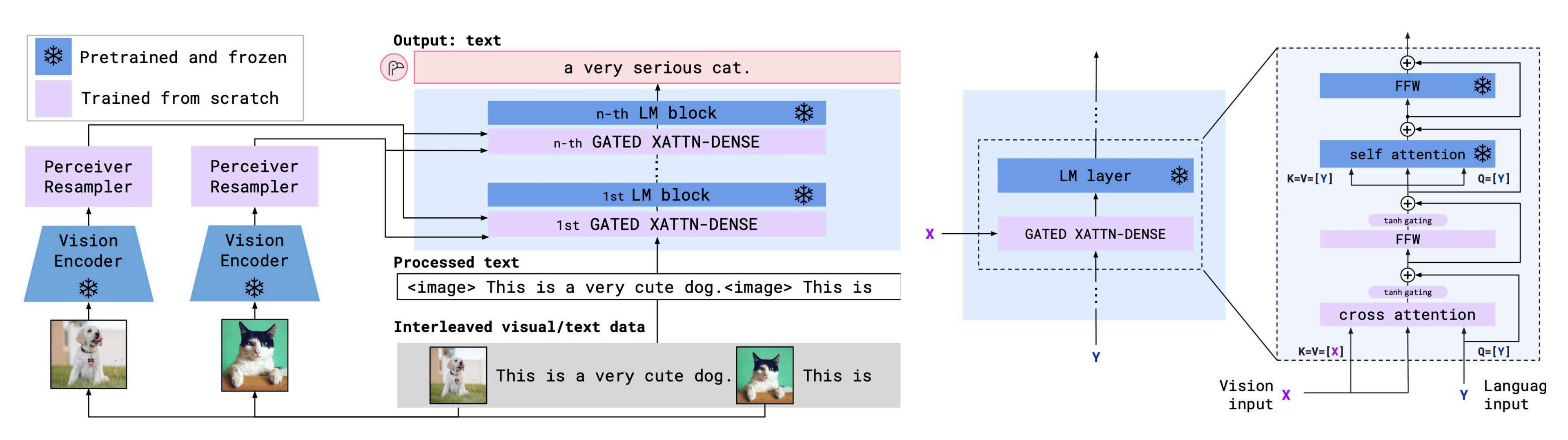
Vision-Language Model (VLM)

- Vision-Language Models (VLMs) are an emerging class of Generative
 Als that aim to understand and generate content that combines both
 visual and textual information. They are designed to capture the intricate
 relationships between visual data, such as images or videos, and
 associated language descriptions or textual data.
- Multimodal Learning: VLMs are at the forefront of multimodal learning, where the goal is to build models that can process and relate information from more than one modality in this case, vision (images, videos) and language (text).
- Joint Embedding Space: They typically work by mapping both visual and textual inputs into a common embedding space where the relationship between the two can be learned and understood.
- Cross-Modal Understanding: VLMs are trained to not only understand each modality independently but also to perform cross-modal tasks, such as image captioning, text-to-image generation, and visual question answering (VQA).



Vision-Language Model (VLM)

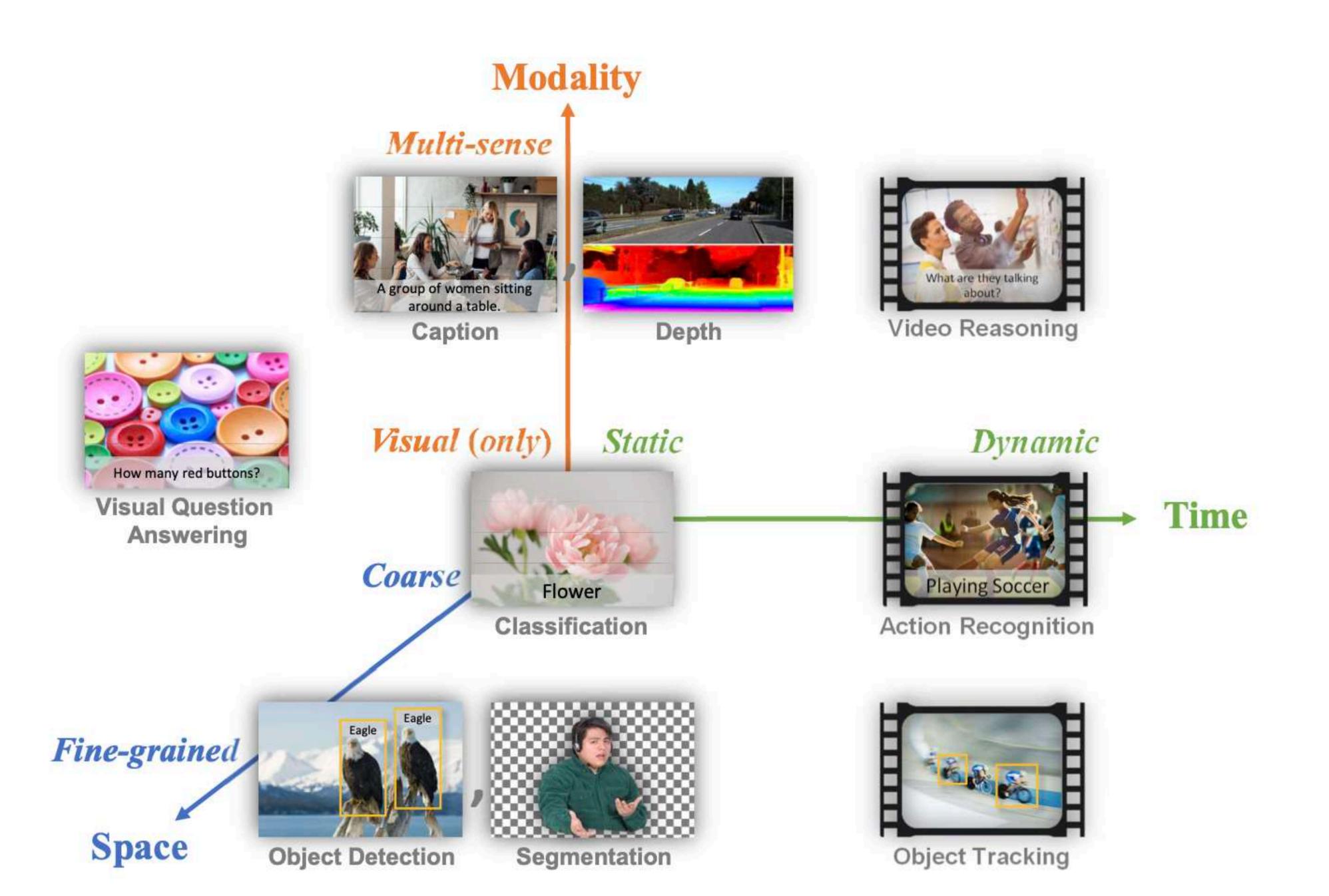
Flamingo: a Visual Language Model for Few-Shot Learning



Flamingo (Alayrac et al., 2022)

Multi-Modalities Vision-Language Model

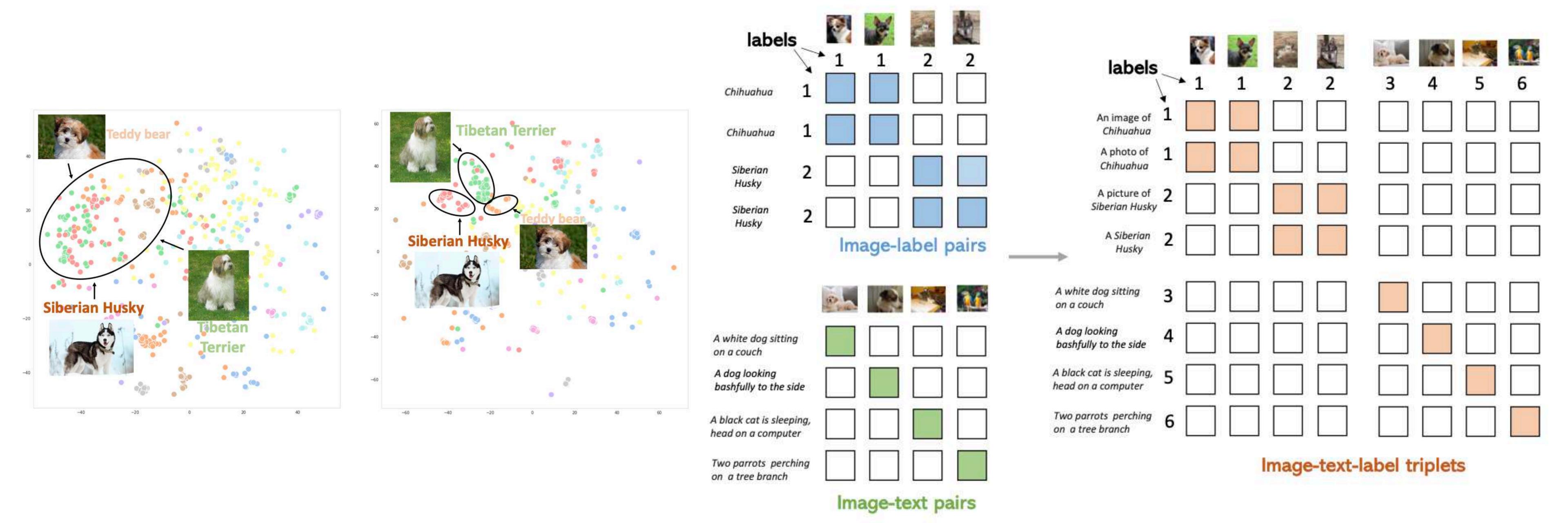
Florence: A New Foundation Model for Computer Vision (Lu et al., 2022)



- STOA Multi-Model Vision-Language foundation model.
- Unified learning
- Transformer architectures and adaption
- ImageNet-1K zero-shot top-1 83.74%
- Release a dataset FLD-900M, containing 900M images and 900M texts.
- Model size efficiency (only 647M)
- Combine multiple previous works:
 - Model design: Dynamic DETR (Dai et al., 2021)
 - UniCL (Yang et al., 2022),
 - ZeRO optimizer, mix-precision training, etc.,

Multi-Modalities Vision-Language Model

Florence: A New Foundation Model for Computer Vision (Lu et al., 2022)



UniCL (Yang et al., 2022)



Introduction

Video Action Anticipation

Introduction to Diffusion Models

Synthesis and Future Directions

Conclusion and Discussion

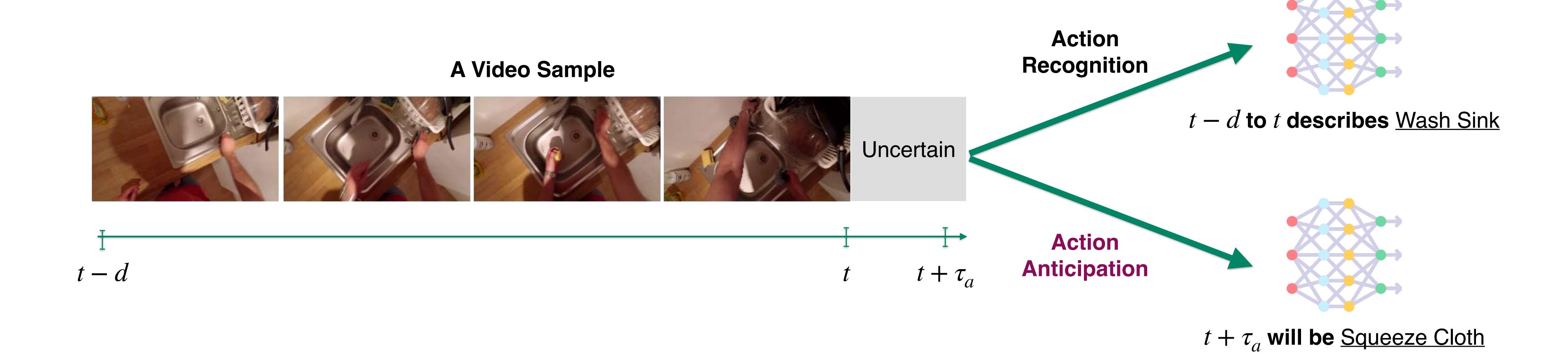
Challenge in Videos

Curse of Dimensionality

- · Video datasets present unique challenges for analysis and modeling, primarily due to their high-dimensional nature.
 - · High Dimensionality: Managing large data volumes without losing crucial information.
 - Temporal Dynamics: Modeling frame sequences and changes over time.
 - Computational Demands: Requires more resources than text data processing.
- Comparing <u>video data</u> with <u>text data</u>:
 - Dimensionality and Encoding
 - Text Data: Encoded as vectors, typically ~20,000 elements for a word, capturing semantic properties.
 - Video Data: Far more complex; a 256 x 256 x 3 (RGB) frame over time T requires ~20,000 x T elements, as each frame is a high-resolution image.
 - Discrete vs. Continuous Nature
 - Text Data: Discrete in nature; words are distinct and ordered units.
 - · Video Data: Continuous in both spatial and temporal aspects, where small changes can significantly alter content.
 - Temporal Sensitivity and Sampling
 - Text Data: Order is key, but lacks direct temporal encoding.
 - Video Data: Highly time-sensitive; frame rates and action durations within frames are critical for interpretation.

Video Action Anticipation

Action Recognition versus Action Anticipation



- The same video, when approached with different problem definitions (i.e., recognition vs. anticipation), can lead to varying objectives.
- This demonstrates that action anticipation requires a different solution than action recognition. Solutions directly adapted from action recognition are often sub-optimal.

What Would You Expect? Anticipating Egocentric Actions with Rolling-Unrolling LSTMs and Modality Attention

A. Furnari, G. M. Farinella

(egocentric action anticipation)

degli STUDI di CATANIA

UNIVERSITÀ



Encoding...

ENCODING

ANTICIPATION

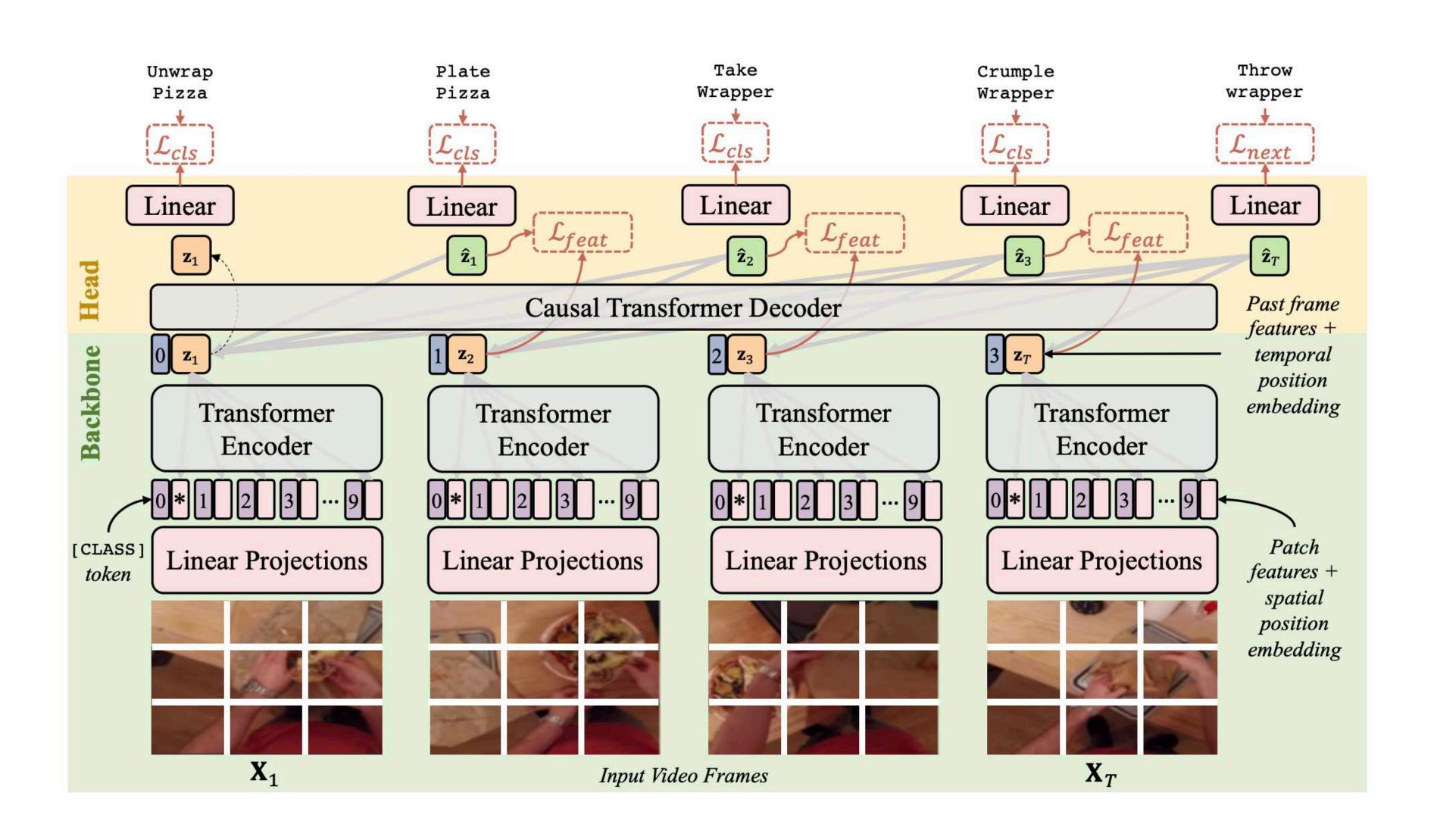
ACTION

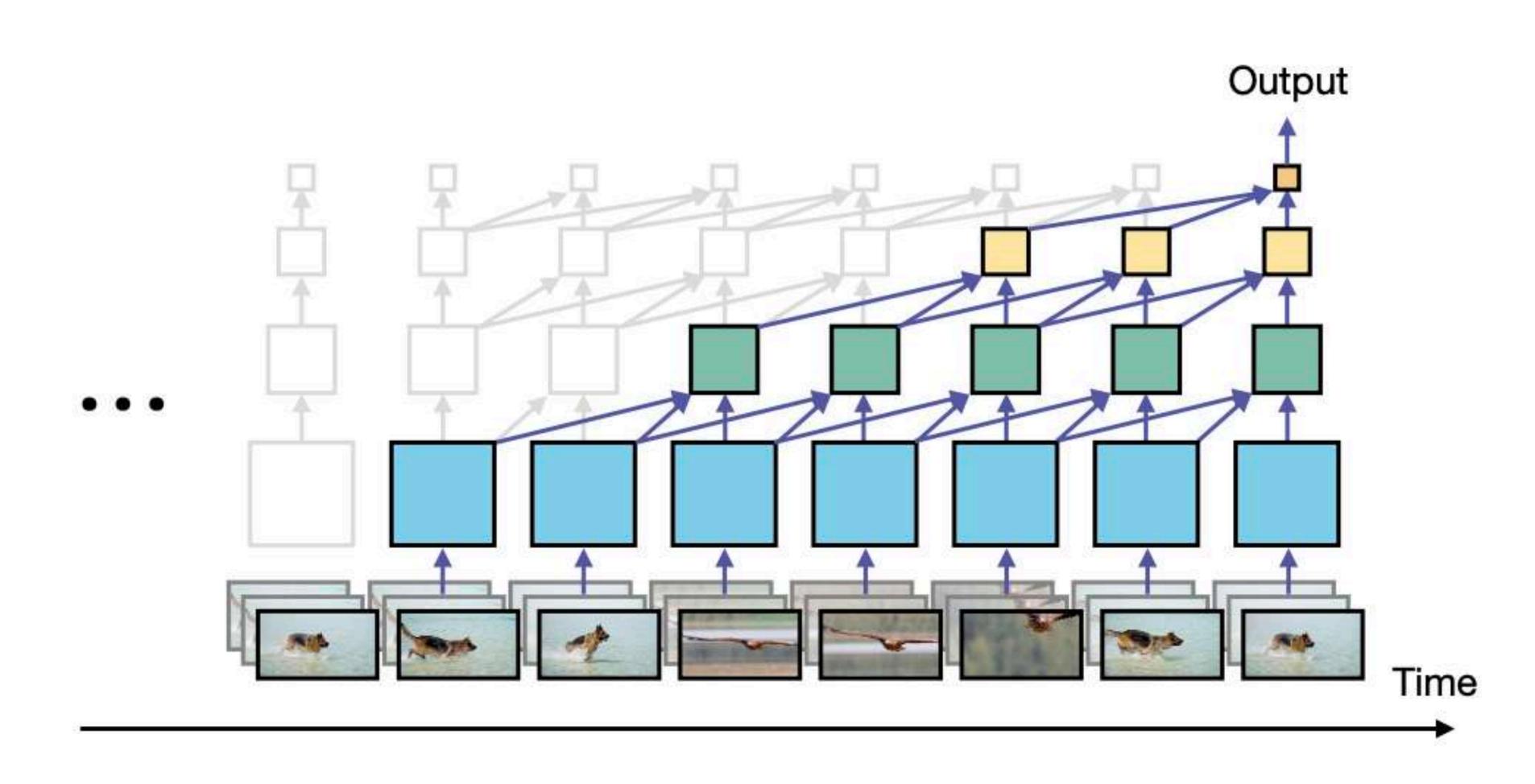
2 1.75 1.5 1.25 1 0.75 0.5 0.25

TAKE SPATULA



Baselines for Video Action Anticipation





AVT (Girdhar et al., 2021)

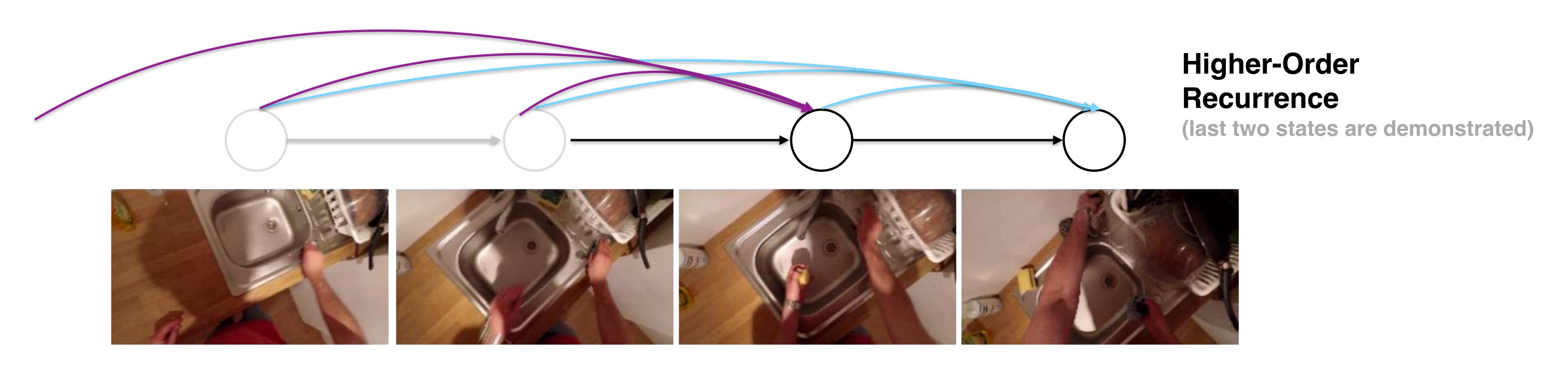
MeMViT (Wu et al., 2022)

 Unlike the previous work, our general strategy for solving this involves modeling long-term dependencies and capturing subtle evidence that is conditional to the targets.

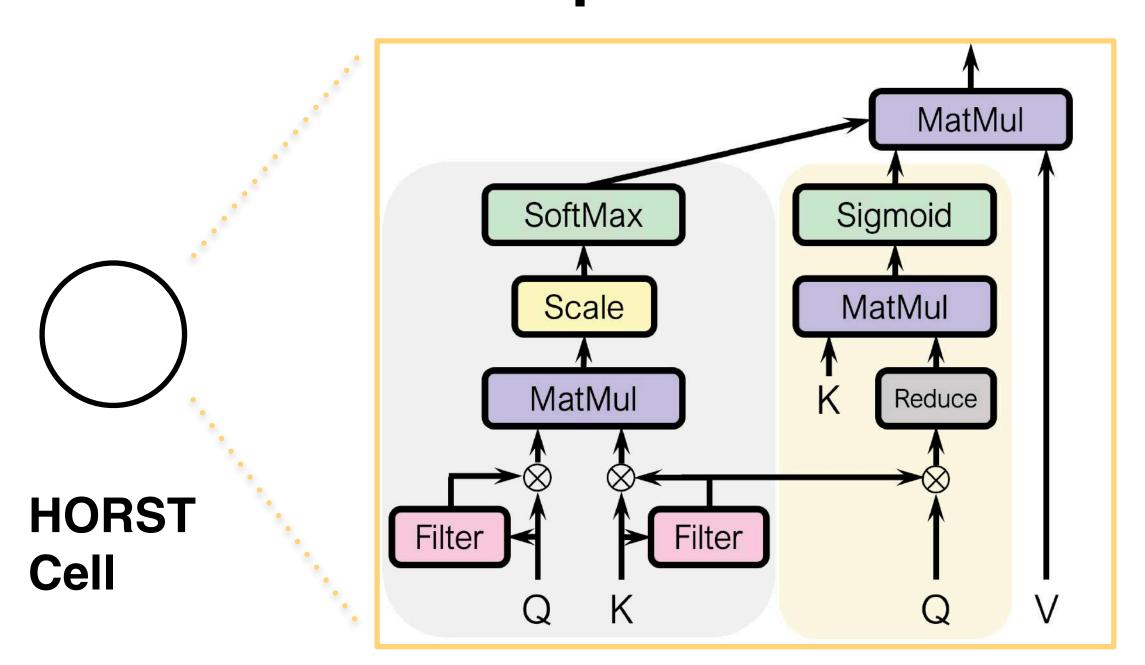
Higher Order Recurrent Spatial-Temporal Transformer (HORST)

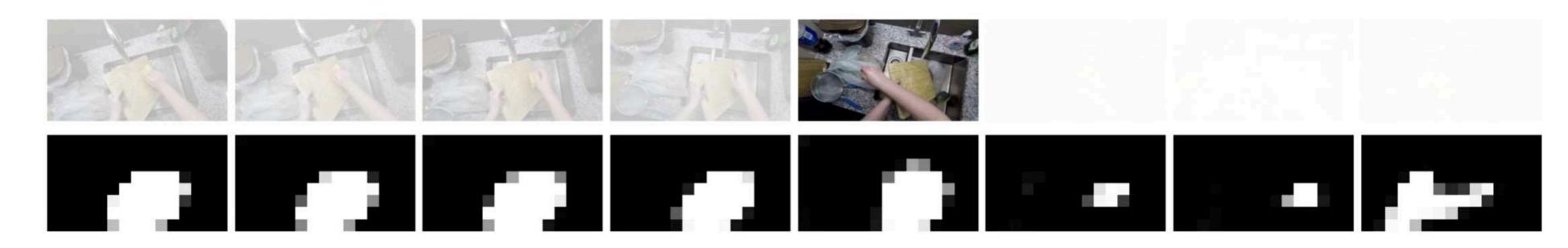
Tai, Tsung-Ming, et al. "Higher-Order Recurrent Network with Space-Time Attention for Video Early Action Recognition." (ICIP 2022).

Tai, Tsung-Ming, et al. "NVIDIA-UNIBZ Submission for EPIC-KITCHENS-100 Action Anticipation Challenge 2022." arXiv preprint arXiv:2206.10869 (2022).



Spatial-Temporal Decomposition Attention





Ground Truth: Wash Board; Top-5 Predictions: Wash Board, Put-Down Board, Turn-Off Tap, Put-Down Sponge, Stir Pan.

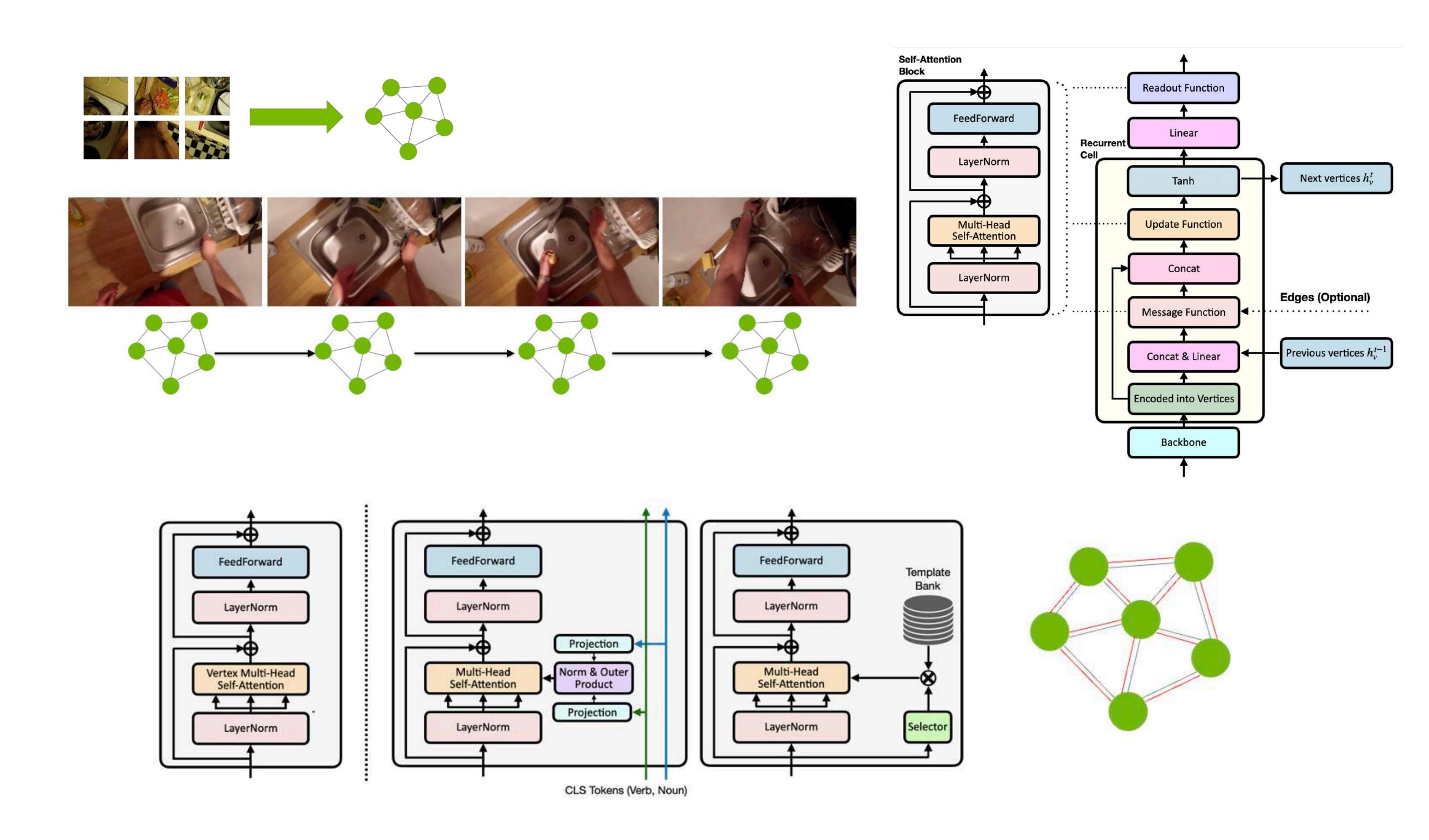
- As N-Gram modeling, HORST leverages higher-order recurrent design. Each timestep, the output is computed depending on multiple previous hidden states, instead of the last one.
- Spatial-Temporal attention learns to query the relevant information from higherorder queue with the decomposition attention of spatial and temporal branches



Unified Recurrent Modeling for Video Action Anticipation

Tai, Tsung-Ming, et al. "Unified recurrence modeling for video action anticipation." (ICPR 2022)
Tai, Tsung-Ming, et al. "NVIDIA-UNIBZ Submission for EPIC-KITCHENS-100 Action Anticipation Challenge 2022." arXiv preprint arXiv:2206.10869 (2022).

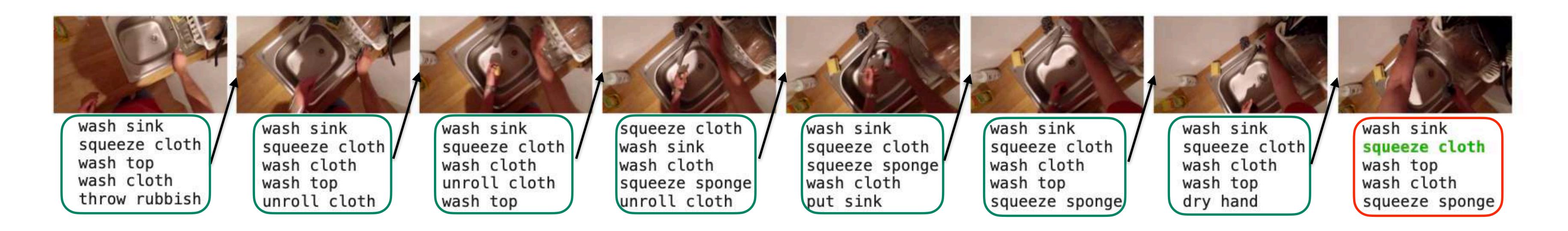
- We proposed MPNNEL, first encodes each frame onto a spatial graph representation, then leverages the message passing framework to learn the temporal propagation.
 - The message-passing framework is composed of Message Function, Update Function, and Readout Function.
 - We implement all of them by the multi-head self-attention.
- The vertex of the spatial graph can be also augmented by our novel edge learning designs, which are
 - Class-Tokens Projection (CTP): which uses two learnable class tokens, supervised from verbs and nouns, and forms the estimated adjacency matrix to augment the edges.
 - Template Bank (TB): Select the learnable templates by the query and using them to augment the edges.
- Using the self-attention here can be treated as learning the information routing between vertices by attention weightings implicitly. Our edge learnings provide a flexible way to further augment the edges explicitly.



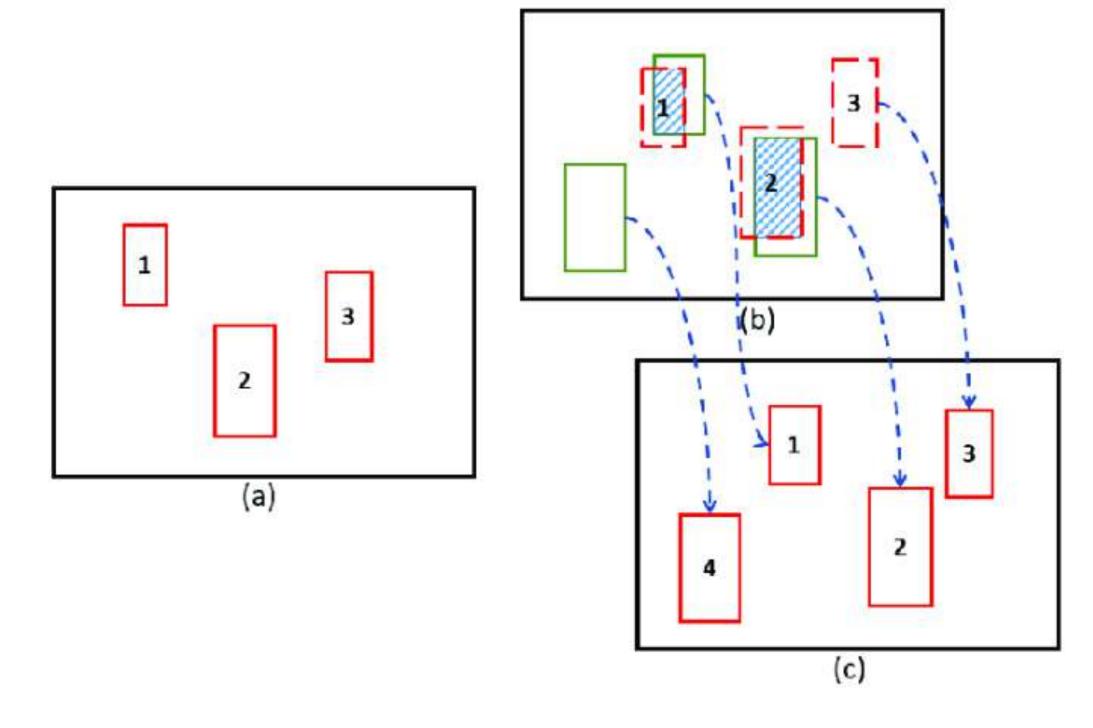


Inductive Attention for Video Action Anticipation

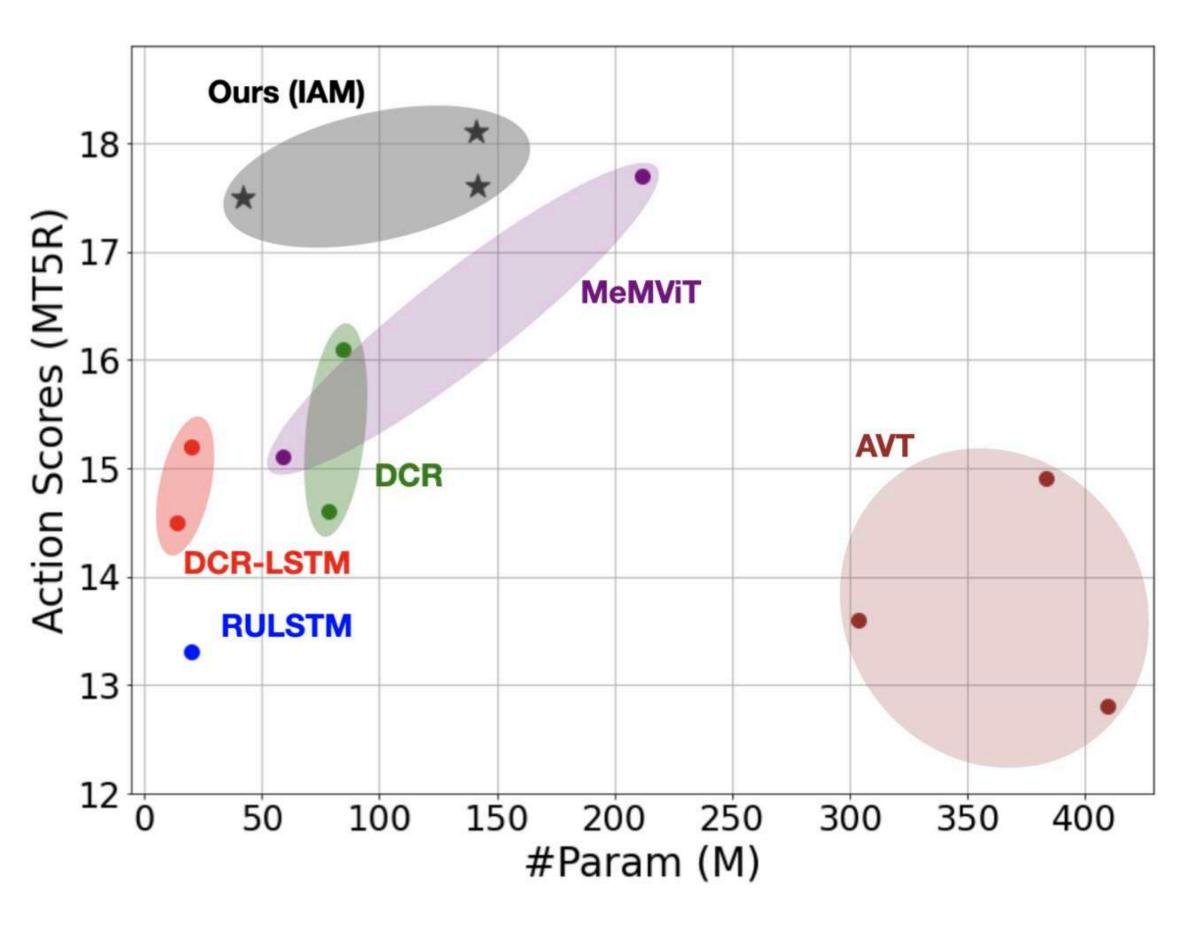
Tai, Tsung-Ming, et al. "Inductive Attention for Video Action Anticipation." arXiv preprint arXiv:2212.08830 (2022).



- We further explore what should be the suitable "query" in attention to drive the higher-order recurrent model.
- In addition, we augment the higher-order recurrent model. Instead of modeling in the latent space, we explicitly utilize the previous predictions as reference (as attention "keys").
- From the experimental results, we found our design can lead to more efficient design in model size and generalize better on the unseen test set.



(For example, object tracking explicitly take the previous prediction as a reference)





Inductive Attention for Video Action Anticipation

(Under Review)

 The first model provides the transparent explainability of previous actions to the future

prediction.

- Strong generalizability on the unseen tests.
- Evaluate on the large-scale EPIC-Kitchens dataset, which contains the 3806 daily kitchen activities.



drink beer pour beer move glass take mat insert cloth

Top-5 Positive

turn-on candle

Top-5 Negative

check basil

fold cloth

take olive

take milk

open drink

remove tofu

flatten box

use cloth



move glass

Top-5 Positive

pour beer

open bottle

close bottle

remove shelf

check cherry

Top-5 Negative

check blueberry

take powder:washing

put bottle

check box

put cap

pour beer put bottle drink beer move glass close bottle

close oil

pour beer

put oil

pour oil

shake bottle

Top-5 Negative

crush tablet

check box

choose egg

peel dough

sort mixture

put bottle fold cloth drink beer close bottle Top-5 Positive

Top-5 Positive close oil close squash take rosemary pour squash close vinegar

pour beer

Top-5 Negative press alarm take salami crush tablet hold bowl take battery

Top-5 Positive close bottle pour oil put bottle pour sauce mix food

pour beer

put bottle

drink beer

pour water

close bottle

Top-5 Negative add vinegar fold yeast move mushroom open vinegar

close bottle pour oil put bottle pour water mix food

put bottle

close bottle

pour beer

pour water

drink water

Top-5 Positive

Top-5 Negative fold yeast insert cellar:salt insert cellar:salt add vinegar open vinegar check bag

close oil close vinegar pour squash pour salad take rosemary

Top-5 Positive

put bottle

pour beer

drink water

pour water

close bottle

Top-5 Negative insert squash move airer move oven carry clothes take backpack

Inductive **Priors** (Respect to highest model prediction e.g., put bottle

Inputs

Anticipated

Actions

Anticipated

Actions

Topmost

Positive/

Negative



take breadstick

move glass sort mat put mat close drawer move cutlery



insert container move glass insert board:chopping move container insert towel



put mat fold cloth put board:chopping take mat move glass



insert container insert towel move glass move container close container



insert container take container put container move container open container



insert container put container take container open container close container

take container close container

insert container move container open container

insert container take container put container dry ladle take tongs

put bottle

pour beer

pour water

close bottle

shake bottle

(GT: take plate)

Topmost Positive/ Negative Inductive

> **Priors** (Respect to highest model prediction e.g., insert container)

Top-5 Positive sort mat fold cloth close drawer insert towel take mat

Top-5 Negative filter lid wash food filter potato lift jar filter container

Top-5 Positive insert container put container put plate open container move container

Top-5 Negative shake peeler:potato close mat close processor:food pour celery hang hand

Top-5 Positive fold cloth sort mat take mat drink beer put board:chopping

Top-5 Negative pat colander put toaster take watch adjust container cut finger: lady

Top-5 Positive insert container move container put container close container open container

Top-5 Negative pour celery shake board:chopping shake peeler:potato search towel feel board:chopping

Top-5 Positive insert containe move container put container take mushroom move cup

Top-5 Negative carry bag pull cloth fold mat pour celery move napkin

Top-5 Positive insert container wash lid insert olive wash container wrap container

Top-5 Negative pour celery gather garlic close plate shake board:chopping scrape garlic

Top-5 Positive close container close olive open container move board:chopping take container

Top-5 Negative adjust container sort kiwi attach tray remove can turn-off processor:food



take cloth

Top-5 Negative

shake straw

throw kale

hang hand

move stand

shake rubbish



take cloth

shake straw

shake caper

wash plug

pat napkin

shake rubbish

Top-5 Negative



squeeze sponge

Top-5 Negative

shake straw

brush sink

adjust straw

shake caper

take alarm





take cloth

Top-5 Negative

shake straw

throw squash

shake rubbish

wash straw

drop sponge



take cloth

Top-5 Negative

shake straw

brush sink

move filter

move coffee

shake rubbish

dry hand

hang hand

move coffee

shake straw

squeeze caper

shake rubbish

Top-5 Negative





Top-5 Positive squeeze cloth	Top-5 Positive wash cloth	Top-5 Positive wash cloth	Top-5 Positive squeeze cloth	Top-5 Positive wash cloth	Top-5 Positive squeeze cloth	Top-5 Positive squeeze cloth
wash sink	wash top	unroll cloth	unroll cloth	wash sink	wash plug	wash cloth
unroll cloth	unroll cloth	dry hand	wash plug	dry hand	wash cloth	hang hand
gather rubbish	dry hand	wash hand	pull cup	take cloth	pull cup	
put sink	wash plug	wash top	hang hand	put cloth	hang hand	pull cup
Put Sink	wash plug	wash top	Hang Hand	put Cloth	nang nand	squeeze sponge
Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative
apply bowl	let-go sponge	gather meat	sort broccoli	throw pear	roll omelette	roll lemon
gather broccoli	mix lettuce	gather broccoli	eat pork	put sheets	wrap box	pat omelette
insert bean:green	eat pork	eat broccoli	coat oil	pat omelette	eat broccoli	eat pork
check pie	check box	wash food	brush dough	unroll omelette	turn ring:onion	turn-on liquid:wash
close freezer	soak plate	insert pan:dust	eat broccoli	move omelette	pat noodle	roll omelette
Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive
wash top	squeeze cloth	squeeze cloth	wash sink	squeeze cloth	wash sink	wash sink
wash cloth	unroll cloth	wash sink	wash cloth	squeeze sponge	dry hand	wash top
throw food	wash oven	wash floor	squeeze sponge	wash top	take cloth	dry hand
take cloth	shake hand	drink beer	wash top	wash oven	wash oven	wash oven
throw rubbish	gather rubbish	remove rubbish	wash sponge	put sponge	gather rubbish	wash hand
Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative	Top-5 Negative
shake rubbish	peel squash	insert banana	shake peach	dry scissors	fill liquid	adjust plate
wash plug	insert nut	peel sausage	squeeze garlic	search salt	hang hand	dry sponge
shake straw	take nut	unroll omelette	carry bag	turn-on water	squeeze can	move funnel
dry cloth	cut apple	cut banana	wash straw	throw squash	shake liquid	empty container
hang hand	look bag	hang hand	squeeze caper	hang cup	shake courgette	filter plate
Ton E Docition	Ton & Docition	Ton E Docition	Ton & Donitivo	Ton E Docition	Ton & Docition	Ton & Docitive
Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive	Top-5 Positive
wash top	wash top	wash top	wash sink	squeeze sponge	wash top	wash top
wash cloth	squeeze sponge	unroll cloth	unroll cloth	squeeze cloth	squeeze sponge	wash cloth
squeeze cloth	wash cloth	wash oven	pour liquid:washing	wash sponge	wash cloth	squeeze sponge
squeeze sponge	wash sponge	dry hand	wash hand	turn-off tap	dry hand	take cloth

dry hand

Top-5 Negative

throw squash

smell candle

let-go sponge

soak plate

shake caper

squeeze cloth wash top wash cloth squeeze sponge

wash sink

wash sink squeeze cloth wash top wash cloth squeeze sponge

wash sink squeeze cloth wash top wash cloth





Top-5 Positive

take pizza

apply olive

apply basil

take basil







Top-5 Positive

put pizza

apply olive

open pizza

put fruit

close bin





(GT: take paper) Rank 28 of 3806

put pizza

take pizza insert pizza cut pizza remove pizza

Top-5 Positive	
apply olive	
put mushroom	
put potato	
throw rubbish	
throw skin	

Top-5 Negative

look bottle

sort fork

open kale

put parsley

insert oven

open oven Top-5 Negative check board:chopping search book rip wrap:plastic insert ring:onion

Top-5 Positive take pizza apply olive apply basil check pizza take basil

Top-5 Negative

look sauce

remove jar

insert holder

lift dishwasher

look can

Top-5 Positive put pizza flatten dough apply olive fold bag put potato Top-5 Negative

turn-on dishwasher

sort fork

look box

put fruit

check drawer

scrub scissors

take basil Top-5 Negative eat bread press coriander take candle open corn look box

Top-5 Positive

put pizza

take pizza

apply basil

apply olive

Top-5 Negative pat spatula insert ring:onion move omelette check pie put pie

Top-5 Negative look can take backpack take pear insert cap pull box

Top-5 Positive

apply olive

remove pizza

put pin:rolling

cut pizza

put food

Top-5 Positive put pizza take pizza

insert pizza close bin open pizza

Top-5 Negative cut com pat garlic press onion check ginger check onion

Top-5 Positive

insert cloth

put pizza insert pizza check pizza open pizza wrap plate

Top-5 Negative

cut com

press onion

remove stalk

apply onion

sort onion

Top-5 Positive insert pizza remove pizza

cut pizza flatten dough open oven

Top-5 Negative

remove cucumber

move pen

take toaster

eat broccoli

put pear

Top-5 Positive flatten dough cut pizza apply dough remove pizza

Top-5 Negative

rip wrap:plastic

divide pizza

unroll clothes

move fridge

dry food

insert pizza flatten dough cut pizza apply olive put fruit

Top-5 Positive

Top-5 Negative dry food look bottle sort broccoli

pour bean:green

Top-5 Positive

remove pizza

insert pizza

take pizza

brush dough

Top-5 Positive take pizza insert pizza cut pizza put sausage remove pizza

Top-5 Negative remove floor throw oil throw butter take pear

close paper

Top-5 Positive put pizza take pizza insert pizza drink beer open pizza

Top-5 Negative cut corn sprinkle tomato close oregano put pear scrape skin

put pizza take pizza insert pizza

cut pizza

remove pizza

Top-5 Positive apply olive

cut pizza remove pizza apply basil take aubergine

Top-5 Negative squeeze garlic pat burger squeeze sausage mix sausage

Top-5 Positive take pizza insert pizza check pizza open pizza open oven

Top-5 Negative

cut corn

press onion

apply onion

squeeze garlic

sort onion

Top-5 Positive put pizza insert pizza

Top-5 Negative

cut corn

press onion

apply onion

remove stalk

sort onion

check pizza cut pizza open pizza

flatten dough insert pizza take pizza fold bag divide dough

Top-5 Negative

search plate

lift foil

cut pork

sort pizza

hold tray

Top-5 Positive

cut pizza apply basil Top-5 Negative cut corn attach tray take toaster

rub sauce

shake tomato

Top-5 Positive put pizza open pizza close bin drink beer remove salad

Top-5 Negative

cut corn

close pepper

put shell:egg

press onion

insert shell:egg

Top-5 Negative check board:chopping unroll clothes move fridge scrape salmon

Top-5 Positive apply olive take pizza cut pizza put fruit put food

insert oven

take pizza insert pizza cut pizza remove pizza

put pizza





Introduction

Video Action Anticipation

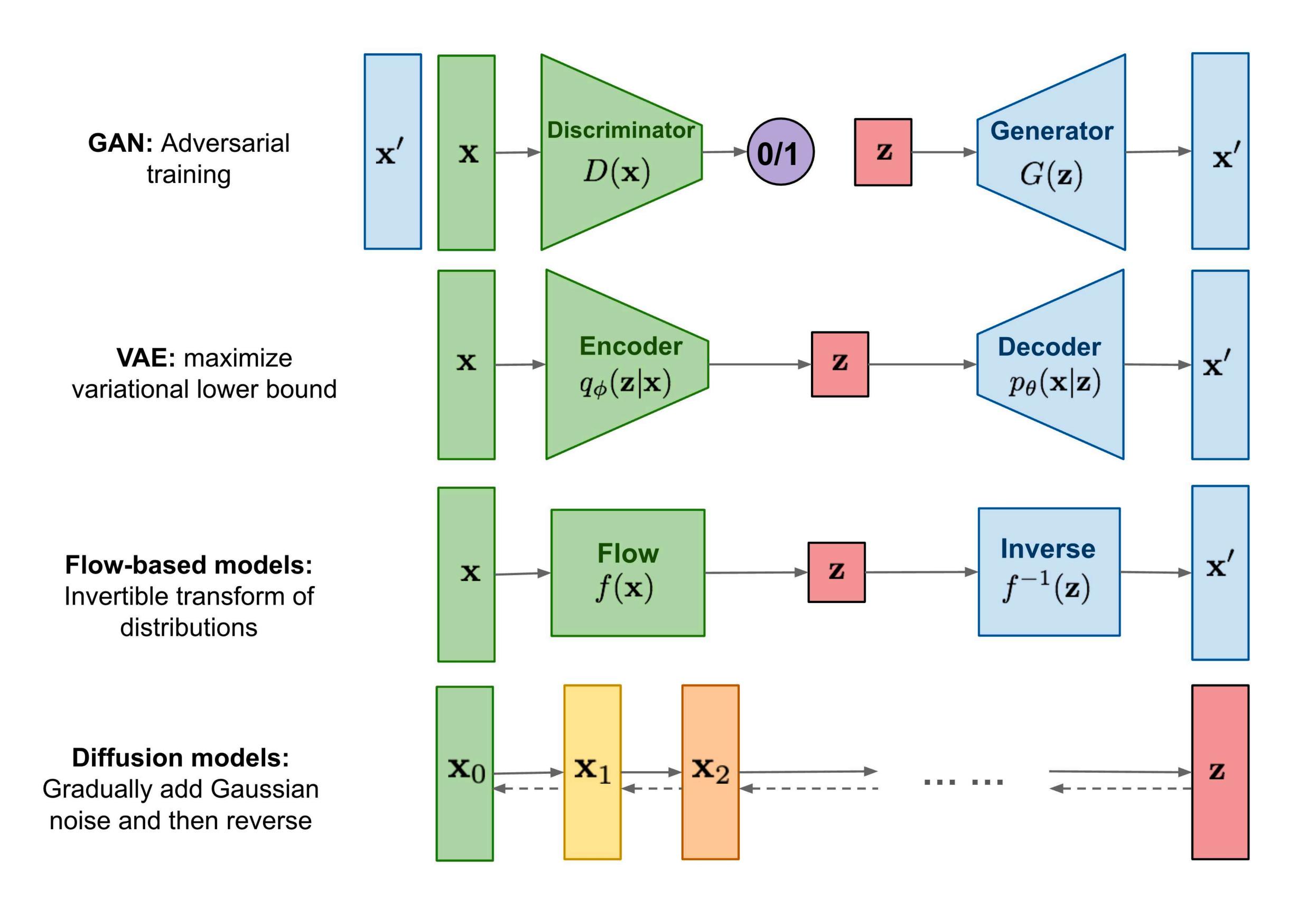
Introduction to Diffusion Models

Synthesis and Future Directions

Conclusion and Discussion

Generative Models

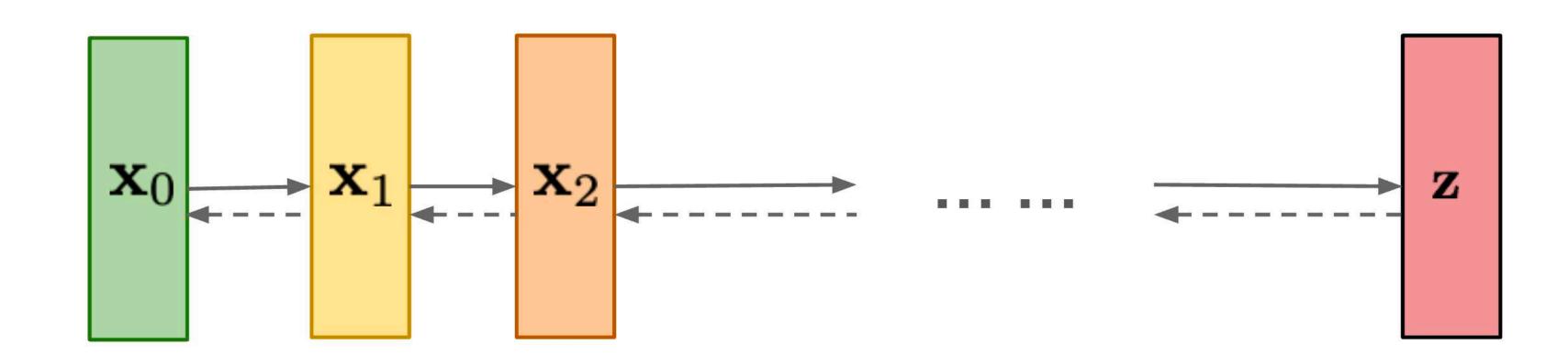
Overview



Markov Chain and Langevin Dynamics

Diffusion models:

Gradually add Gaussian noise and then reverse



Markov chain assumption: $P(x_{t+1} | x_{0:t}) = P(x_{t+1} | x_t)$, which simplifies the way to model the event in continuous time.

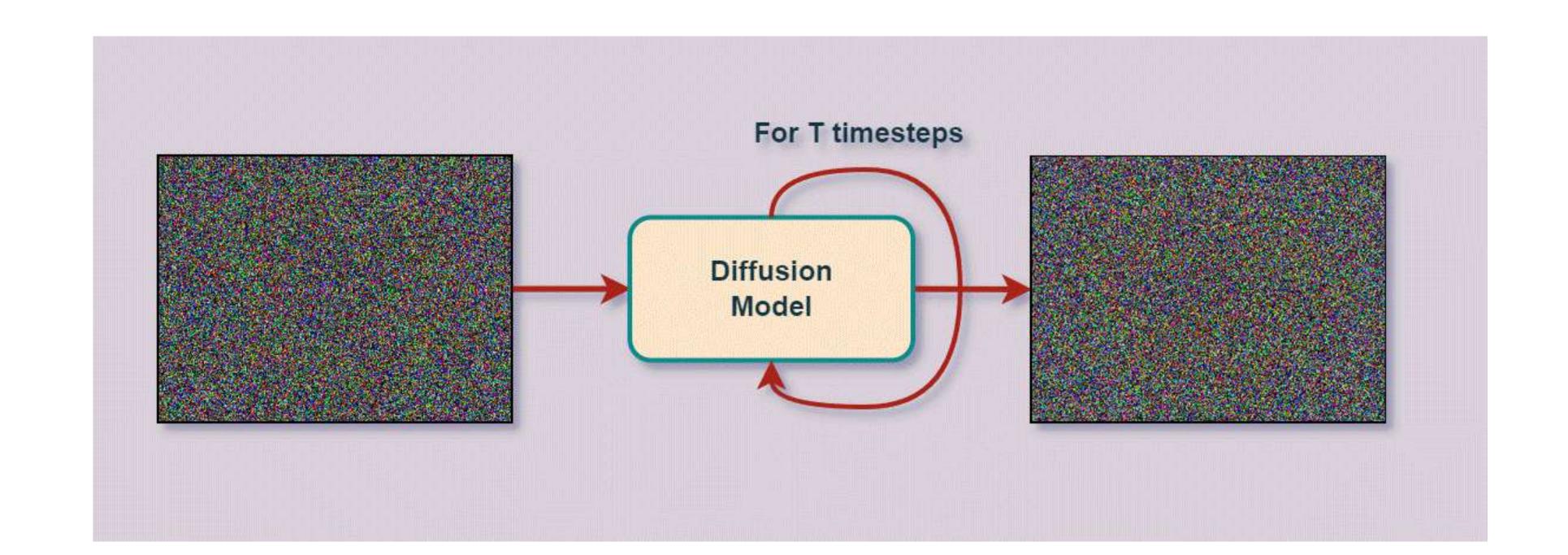
- The diffusion model is based on Langevin (stochastic) Dynamics:
 - $dx_t = \nabla \log p(x_t)dt + \sqrt{2}dB_t$
 - Where B_t is the Brownian Motion with zero mean.
- Utilize the Langevin Dynamics to infer the next small step t + dt:
 - $x_{t+dt} = x_t + \nabla \log p(x) dt + \sqrt{2dt} \cdot N(0,1)$ // this $\sqrt{2dt}$ ensures the following Gaussian noise sampling follows the Brown Motion.
 - $x_{t+dt} = x_t + \nabla \log p(x)dt + N(0,2dt)$
- Here, we can see two different perspectives to model this equation:
 - Continuous Representation: solving the SDE (and its reverse version) and using a neural network to approximate $\nabla \log p(x)$
 - Discrete Representation: solving the Markov chain as each timestep models the small change dt, which is $x_{t+1} = P(x_{t+1} \mid x_t)$

Today's focus!

Denoising Diffusion Probabilistic Models

DDPM (Ho et al., 2020)

- DDPM is a generative model, modeling as a Markov chain, composed of **Diffusion** and **Denoising** phases.
 - Diffusion fused the images with noise.
 - Denoising remove the noise in the image.
- The Markov chain processes for T steps, where T controls the noise strengths in diffusion/denoising phase; T should be sufficiently large (in other words, noise step should be sufficiently small).
- For each step in T,
 - Diffusion utilizes a predefined noises in mixtures.
 - **Denoising** employs a parameterized neural network θ to restore to restore the noisy input.



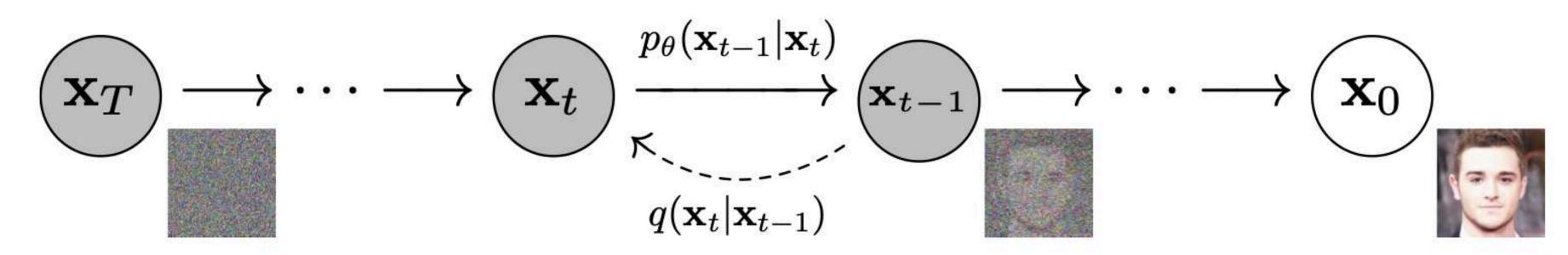
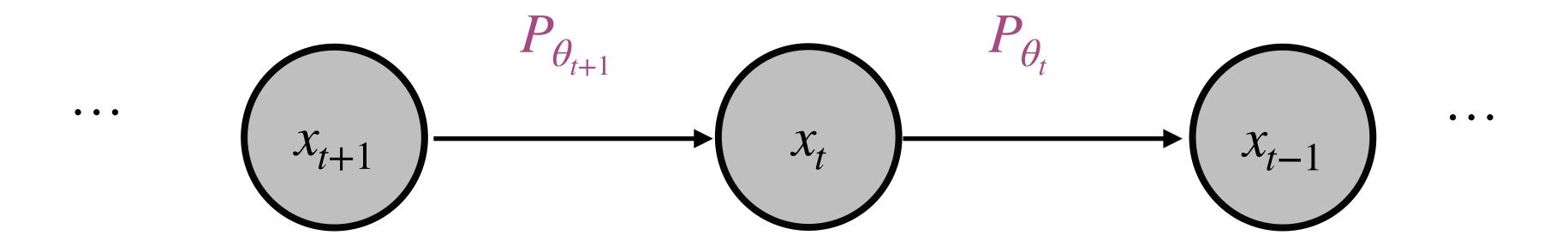


Figure 2: The directed graphical model considered in this work.

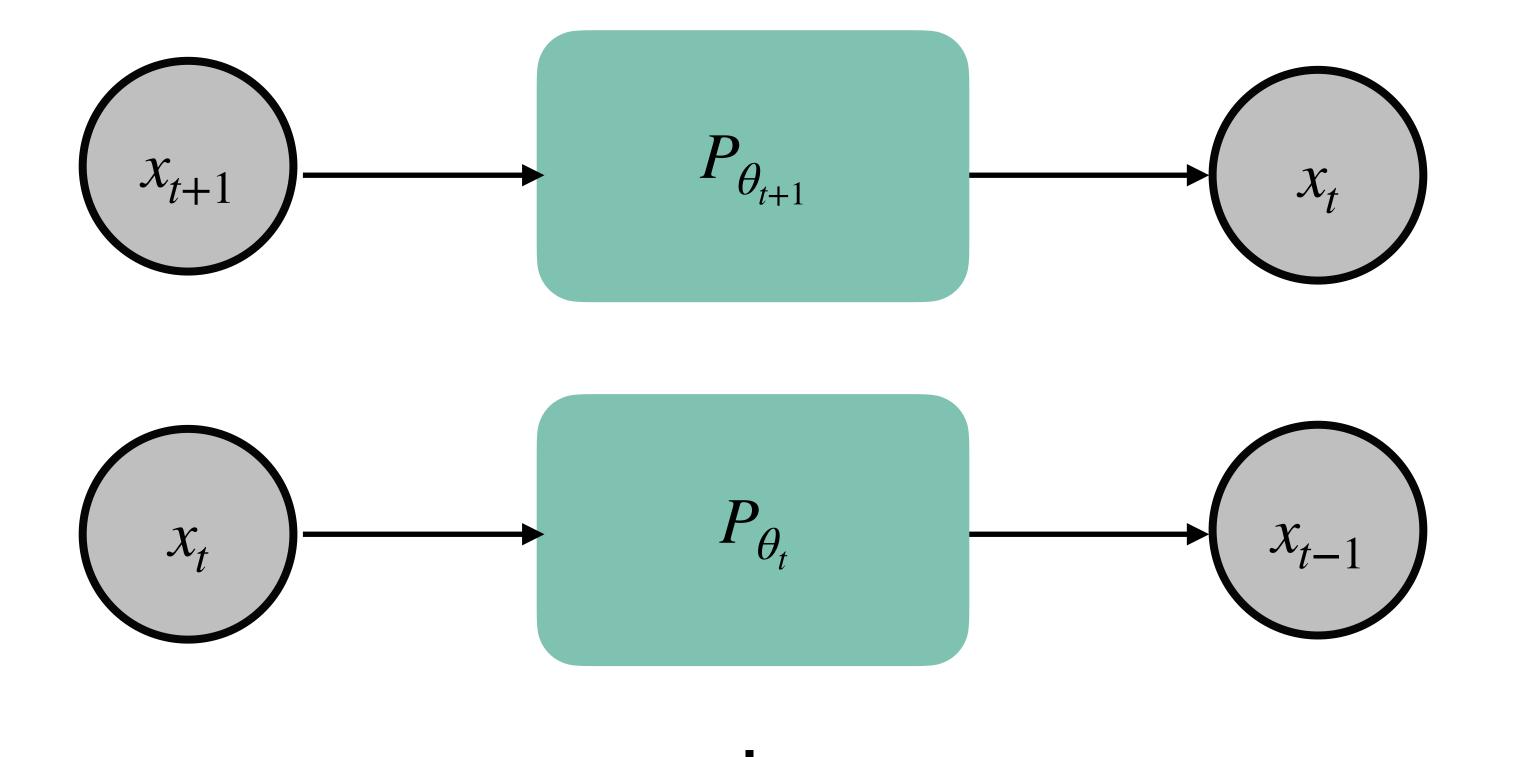
Denoising Diffusion Probabilistic Models

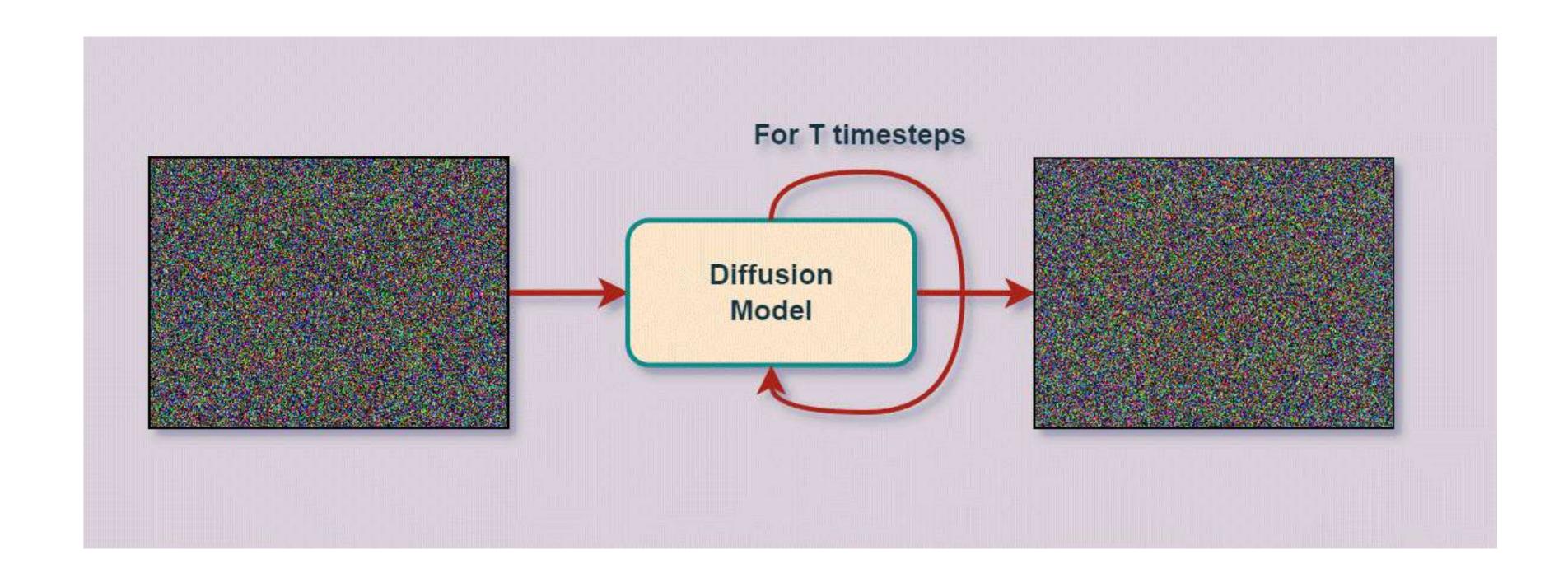
DDPM (Ho et al., 2020)



Note the noise mixtures differently in each diffusion step, As a result, $P_{\theta_{t+1}}$ needs to be different than P_{θ_t}

It is unrealistic to define a total T models for DDPM 9





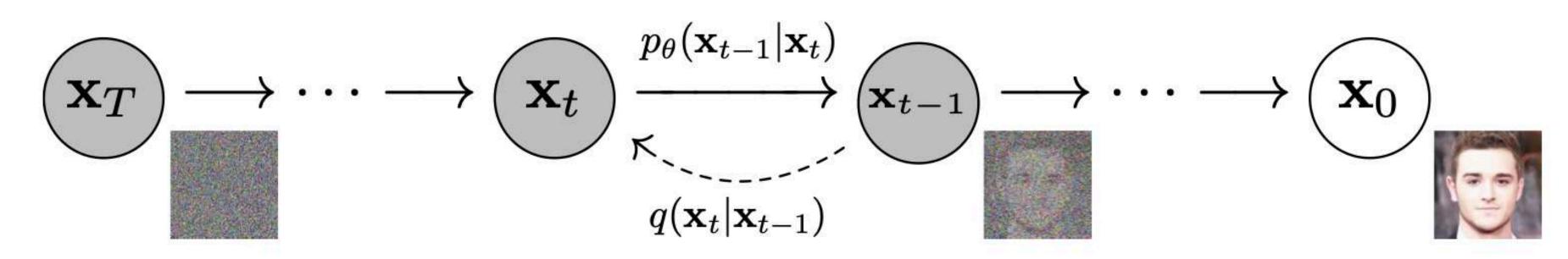
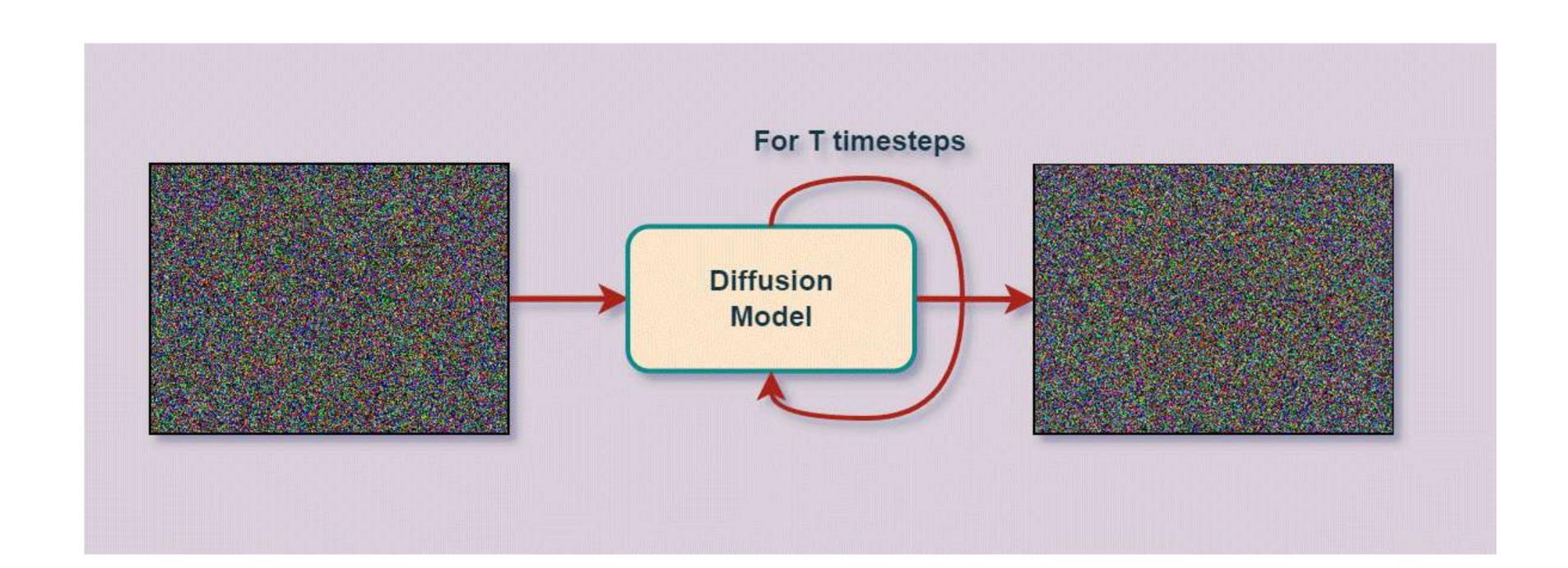


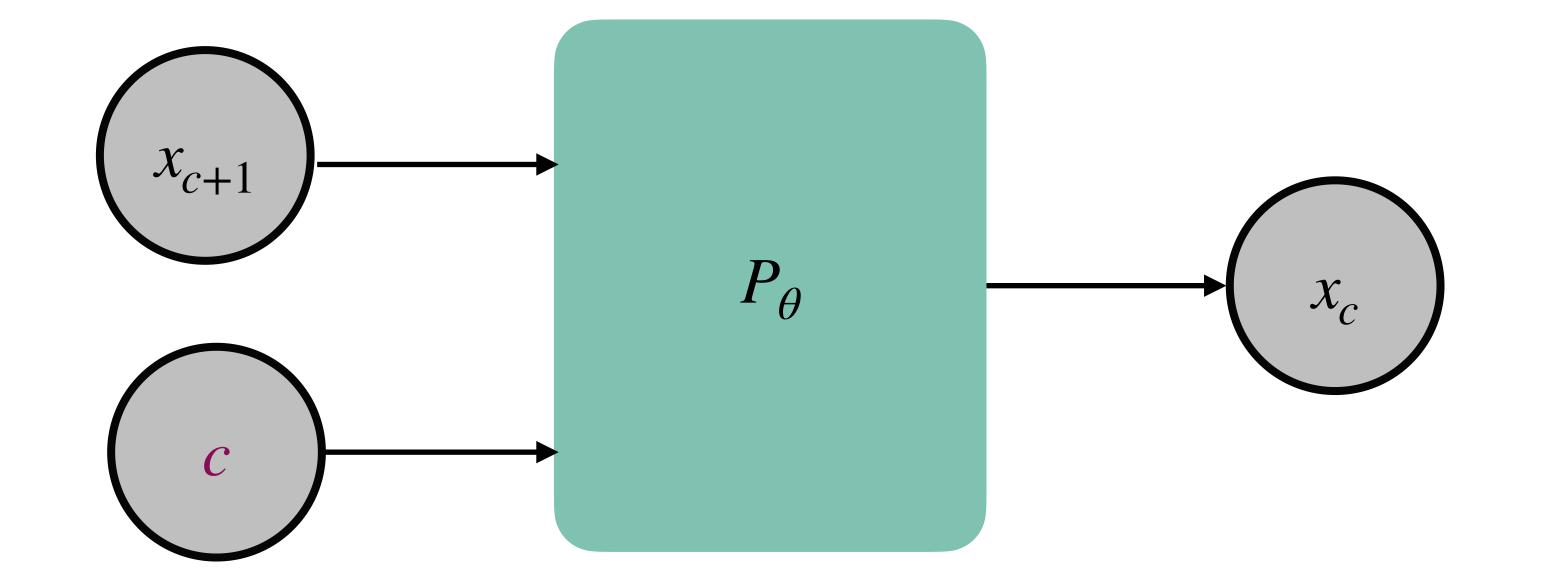
Figure 2: The directed graphical model considered in this work.

Denoising Diffusion Probabilistic Models

DDPM (Ho et al., 2020)

We parameterize the timestep as an auxiliary input (condition) to marginalize the model





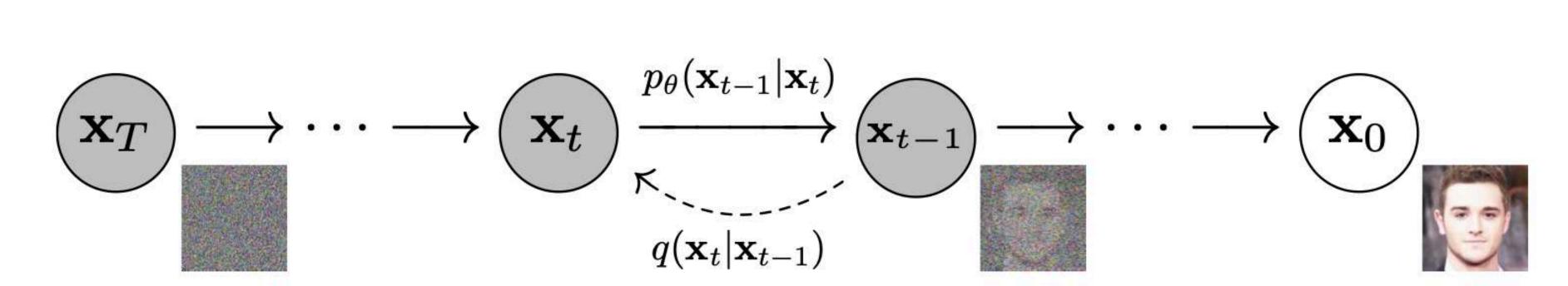


Figure 2: The directed graphical model considered in this work.

Denoising Diffusion Probabilistic Models -- Diffusion Phase

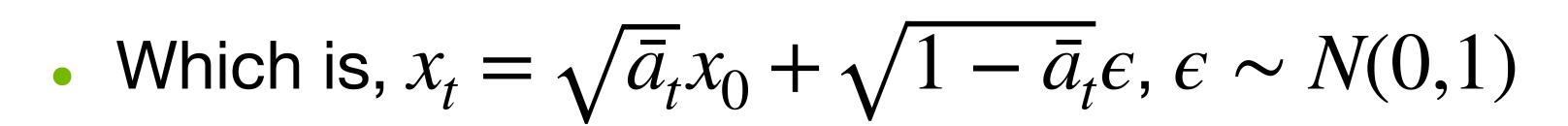
DDPM (Ho et al., 2020)

• The diffusion phase $q(x_{1:T}|x_0) := \Pi_{t=1}^T q(x_t|x_{t-1})$, where $q(x_t|x_{t-1}) := N(\sqrt{\frac{a_t}{a_{t-1}}} x_{t-1}, (1 - \frac{a_t}{a_{t-1}})I)$.

For each diffusion step, the noise mixtures in different strengths rather than the same gaussian noise.

The scheduler is predefined by a <u>decreasing sequence</u> $a_{1:T} = [0,1)^T$. A popular choice is <u>cosine decay</u> (Nichol, Alexander Quinn, and Prafulla Dhariwal. 2021.)

• Since the noise distribution is predefined and known, we can directly approximate $q(x_t|x_0) = N(\sqrt{\bar{a}_t}x_0, \sigma^2 I)$, where $\bar{a}_t = \Pi_{s=1}^t a_s$.



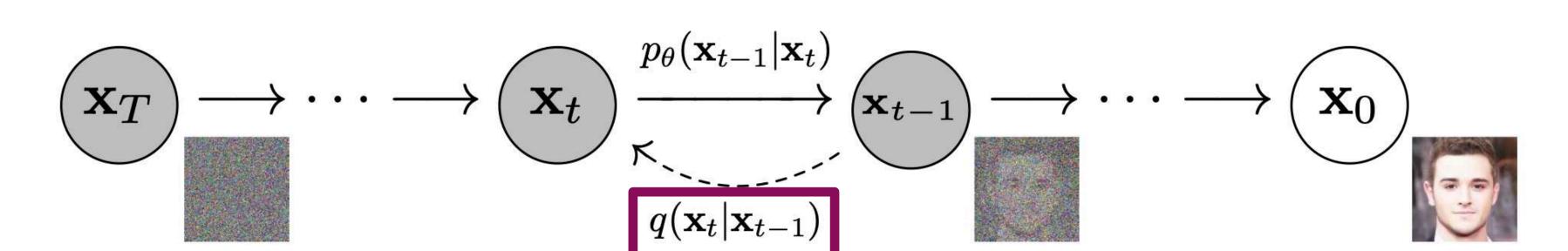




Figure 3. Latent samples from linear (top) and cosine (bottom) schedules respectively at linearly spaced values of t from 0 to T. The latents in the last quarter of the linear schedule are almost purely noise, whereas the cosine schedule adds noise more slowly

Denoising Diffusion Probabilistic Models -- Diffusion Phase

DDPM (Ho et al., 2020)

• The diffusion phase $q(x_{1:T}|x_0) := \prod_{t=1}^T q(x_t|x_{t-1})$, where $q(x_t|x_{t-1}) := N(\sqrt{\frac{a_t}{a_{t-1}}}x_{t-1}, (1-\frac{a_t}{a_{t-1}})I)$.

For each diffusion step, the noise mixtures in different strengths rather than the same gaussian noise.

The scheduler is predefined by a decreasing sequence $a_{1:T} = (0,1]^T$. A popular choice is <u>cosine decay</u> (Nichol, Alexander Quinn, and Prafulla Dhariwal. 2021.)

 Since the noise distribution is predefined and known, we can directly approximate $q(x_t | x_0) = N(\sqrt{\bar{a}_t}x_0, \sigma^2 I)$, where $\bar{a}_t = \Pi_{s=1}^t a_s$.

• Which is,
$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon$$
, $\epsilon \sim N(0,1)$

Present the expected ratio of x_t to x_{t-1} is . , and the variance, σ^2 , is $(1 - \frac{a_t}{\sigma})$.

$$\begin{aligned} x_t &= \sqrt{a_t} x_{t-1} + \sqrt{1-a_t} \epsilon \\ &= \sqrt{a_t} a_{t-1} x_{t-2} + \sqrt{(1-a_t) + a_t} (1-a_{t-1}) \\ &\cdots \\ &= \sqrt{\bar{a}_t} x_0 + \sqrt{1-\bar{a}_t} \epsilon \\ &\text{where } \bar{a}_t = \Pi_{s=1}^t a_s \end{aligned}$$
 Given that,
$$A = N(0)$$

 $A = N(0, \sigma_a^2 I)$ $B = N(0, \sigma_b^2 I)$ Z = A + B $= N(0, (\sigma_a^2 + \sigma_b^2) I)$ And, $N(0, \sigma^2 I) = \sigma N(0, I)$

$$N(0,\sigma^2I) = \sigma N(0,I)$$

We can see when t=T, if and only if $a_T\to 0$, the x_T converges to N(0,1) for all x_0



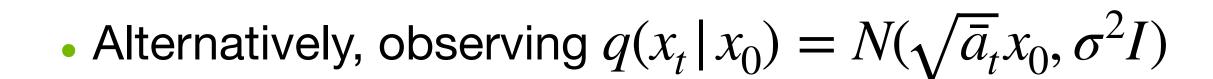
Denoising Diffusion Probabilistic Models -- Denoising Phase

DDPM (Ho et al., 2020)

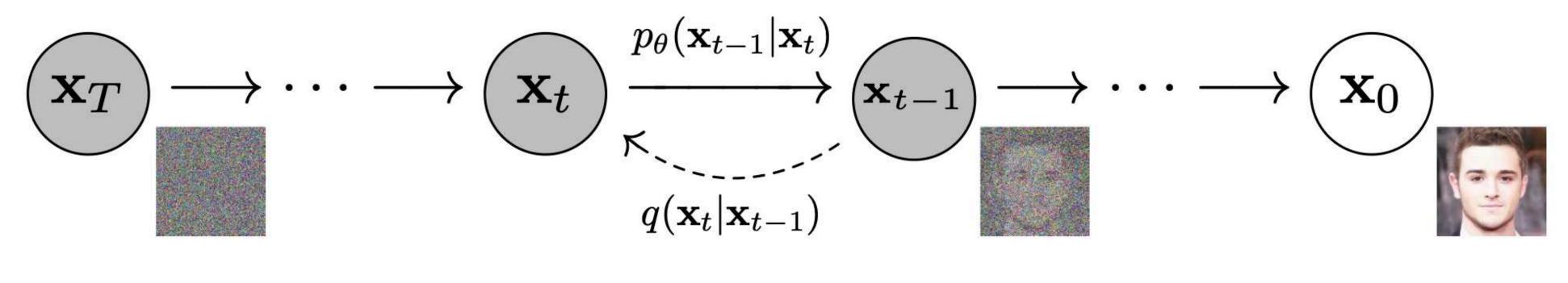
• Reverse process iteratively inference the $p_{\theta}(x_{t-1} | x_t)$ for t = T, ..., 1, parameterized by trainable weights θ .

In DDPM, the reverse chain is defined by
$$x_{t-1} = \frac{1}{\sqrt{a_t}}(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{a}_t}}\epsilon)$$

• Since noise is intractable by only the noisy observation x_t , a neural network $f_{\theta}(x_t,t)$ is deployed to predict the noise ϵ with the noisy observation x_t input.



- We can directly predict $\hat{x_0} = f'_{\theta}(x_t, t)$, and obtain x_t by $x_t = q(x_{t-1} | f'_{\theta}(\hat{x_0} | x_t, t) = N(\sqrt{\bar{a}_t}\hat{x_0}, \sigma^2 I)$
- where $\sigma^2 = \frac{1 a_{t-1}^-}{1 \bar{a}_t} (1 a_t)$ (Sec 3.2 in DDPM paper)



1.
$$x_{t-1} = \frac{1}{\sqrt{a_t}} (x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{a}_t}} f_{\theta}(x_t, t))$$

$$\dots$$

$$x_{t-1} = \frac{1}{\sqrt{a_t}} (x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{a}_t}} f_{\theta}(x_t, t))$$

$$\dots$$

$$\hat{x}_0 = f_{\theta}(x_t, t)$$

$$\vdots$$

$$x_t$$

$$x_{t-1}$$

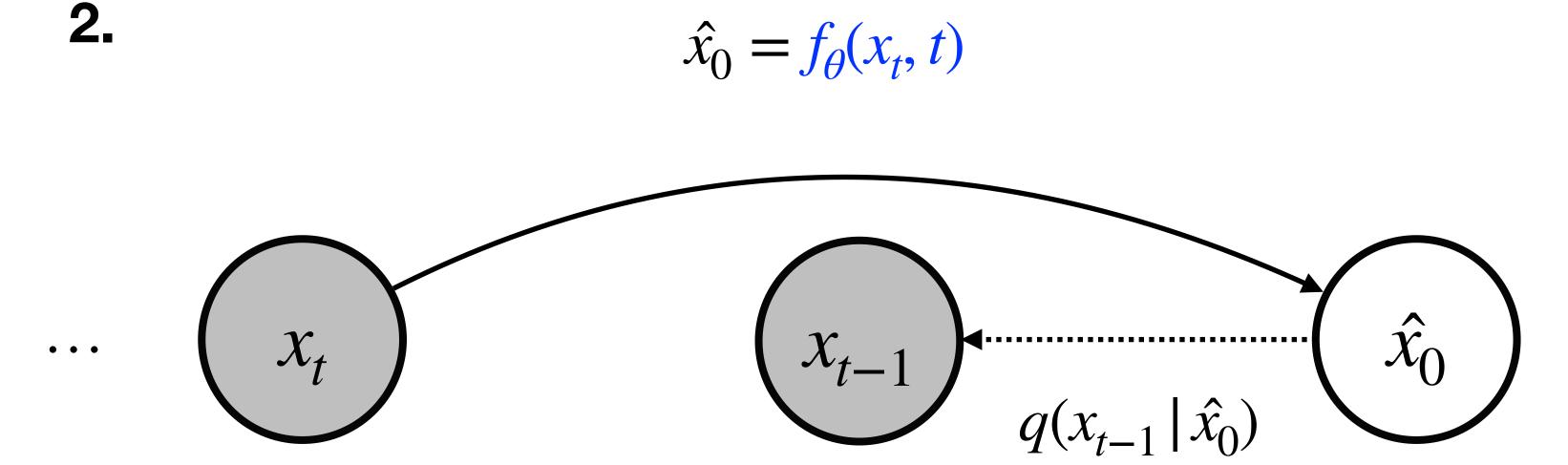
$$x_{t-1}$$

$$\hat{x}_0$$

Denoising Diffusion Probabilistic Models -- Train Stage

DDPM (Ho et al., 2020)

- Method 1, $f_{\theta}(x_t, t)$ is the neural network aims to predict the gaussian noises ϵ by the noisy observation x_t ,
 - Since we can access the ϵ during training, we can supervised f_{θ} by minimizing $||f_{\theta}(\sqrt{\bar{a}_t}x_0 + \sqrt{1-\bar{a}_t}\epsilon,t) \epsilon||_2^2$
- Method 2, $f_{\theta}(x_t, t)$ directly predicts the x_0 , the supervised target, in this case, is straightforward, minimizing $||f_{\theta}'(x_t, t) x_0||_2^2$



Denoising Diffusion Probabilistic Models -- Train Stage

DDPM (Ho et al., 2020)

- Method 1, $f_{\theta}(x_t, t)$ is the neural network aims to predict the gaussian noises ϵ by the noisy observation x_t ,
 - Since we can access the ϵ during training, we can supervised f_{θ} by minimizing $||f_{\theta}(\sqrt{\bar{a}_t}x_0 + \sqrt{1-\bar{a}_t}\epsilon,t) \epsilon||_2^2$
- Method 2, $f_{\theta}(x_t, t)$ directly predicts the x_0 , the supervised target, in this case, is straightforward, minimizing $||f_{\theta}'(x_t, t) x_0||_2^2$
 - x_0 (original image) usually is more stable in numeric computation than ϵ (noise). Method 2 could result in more stable training.

 $\hat{x}_0 = f_{\theta}(x_t, t)$ $\dots \qquad \hat{x}_t \qquad \hat{x}_{t-1} \qquad \hat{x}_0$

High-Level Interpretation of DDPM

"Denoising" Perspective

- From an image with pure noises, DDPM gradually removes the noise bit by bit in the image and reveals the underneath content.
- In this perspective, the neural network heta serves as a Gaussian denoiser.

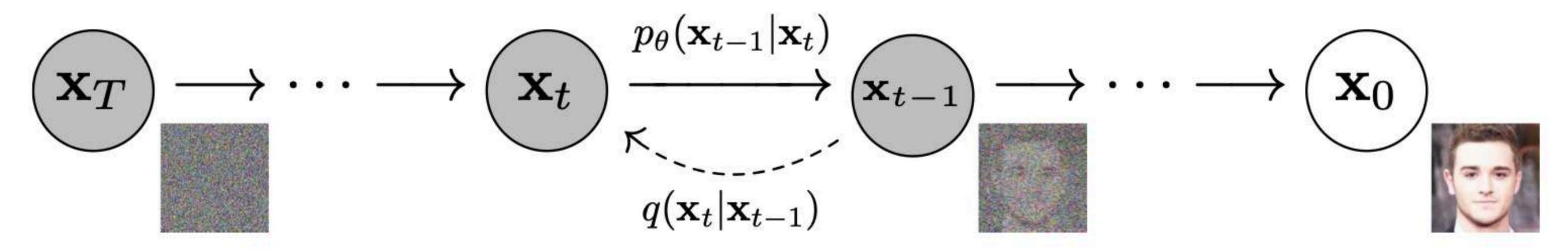


Figure 2: The directed graphical model considered in this work.

High-Level Interpretation of DDPM

"Entropy" Perspective

- Starting from an uncertain state, DDPM gradually excludes the least possible candidate until the only possible state remains.
- ullet In this perspective, the neural network heta learns to reduce the entropy of the corresponding input.
 - (ref: Shannon entropy of a random variable x is defined as E[-logp(x)])

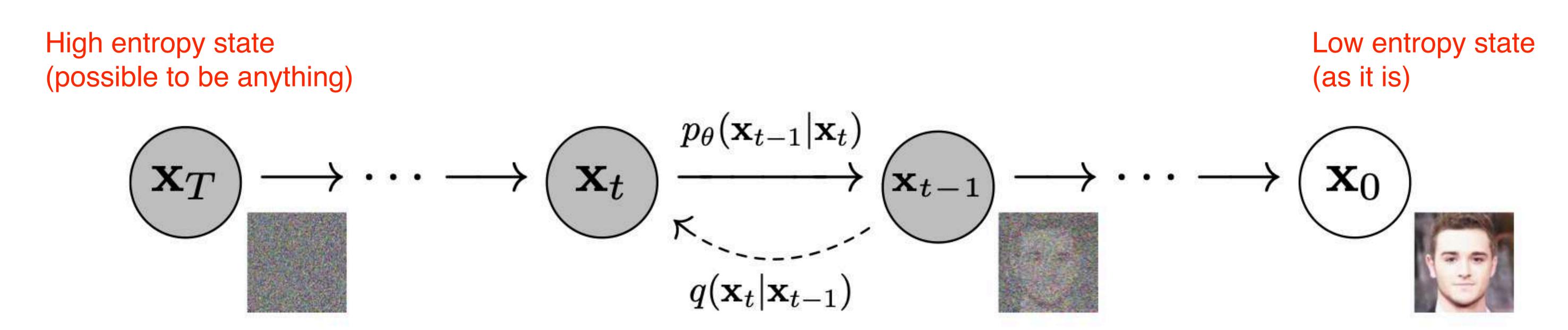
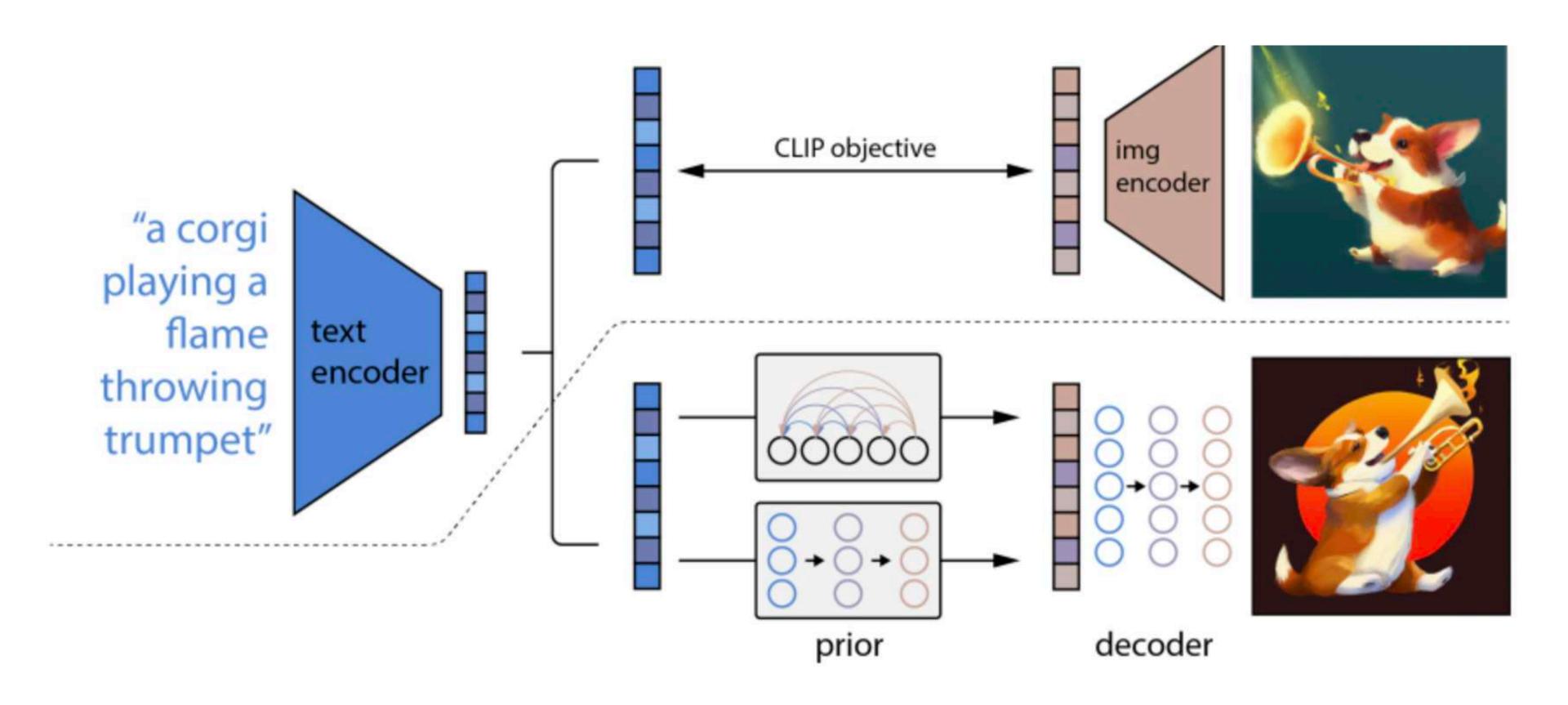


Figure 2: The directed graphical model considered in this work.



Latent Space Conditioning Semantic **Diffusion Process** Map Denoising U-Net $\epsilon_{ heta}$ $|z_T|$ Repres entations Images Pixel Space $au_{ heta}$ crossattention denoising step switch skip connection concat

Stable-Diffusion



LSGM: Score-based Generative Modeling in Latent Space



Introduction

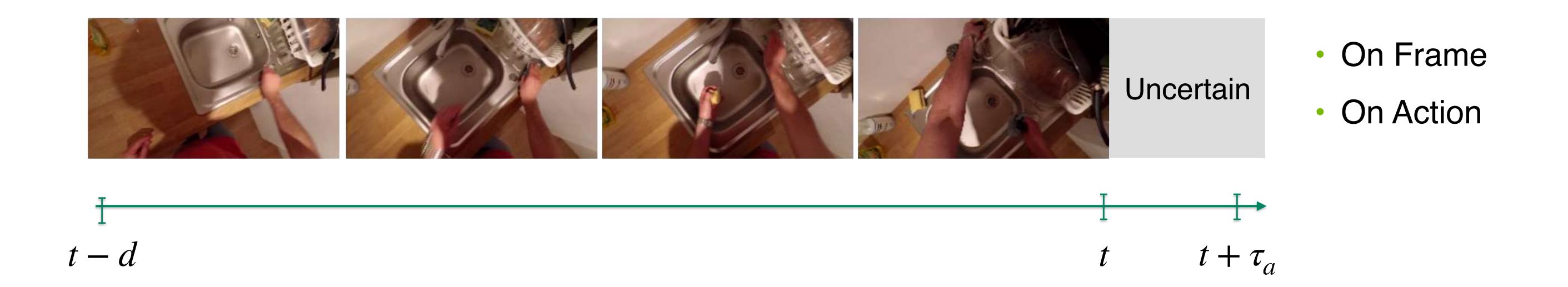
Video Action Anticipation

Introduction to Diffusion Models

Synthesis and Future Directions

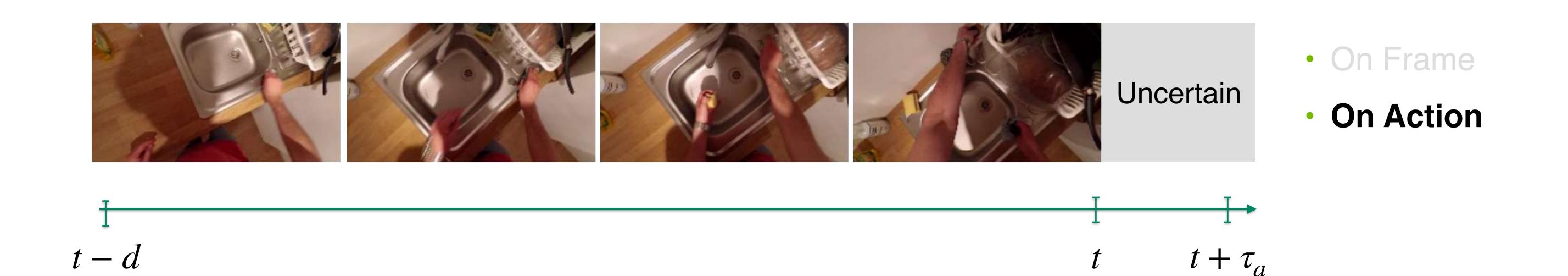
Conclusion and Discussion

Integrating Video Action Anticipation with Diffusion Models





Integrating Video Action Anticipation with Diffusion Models

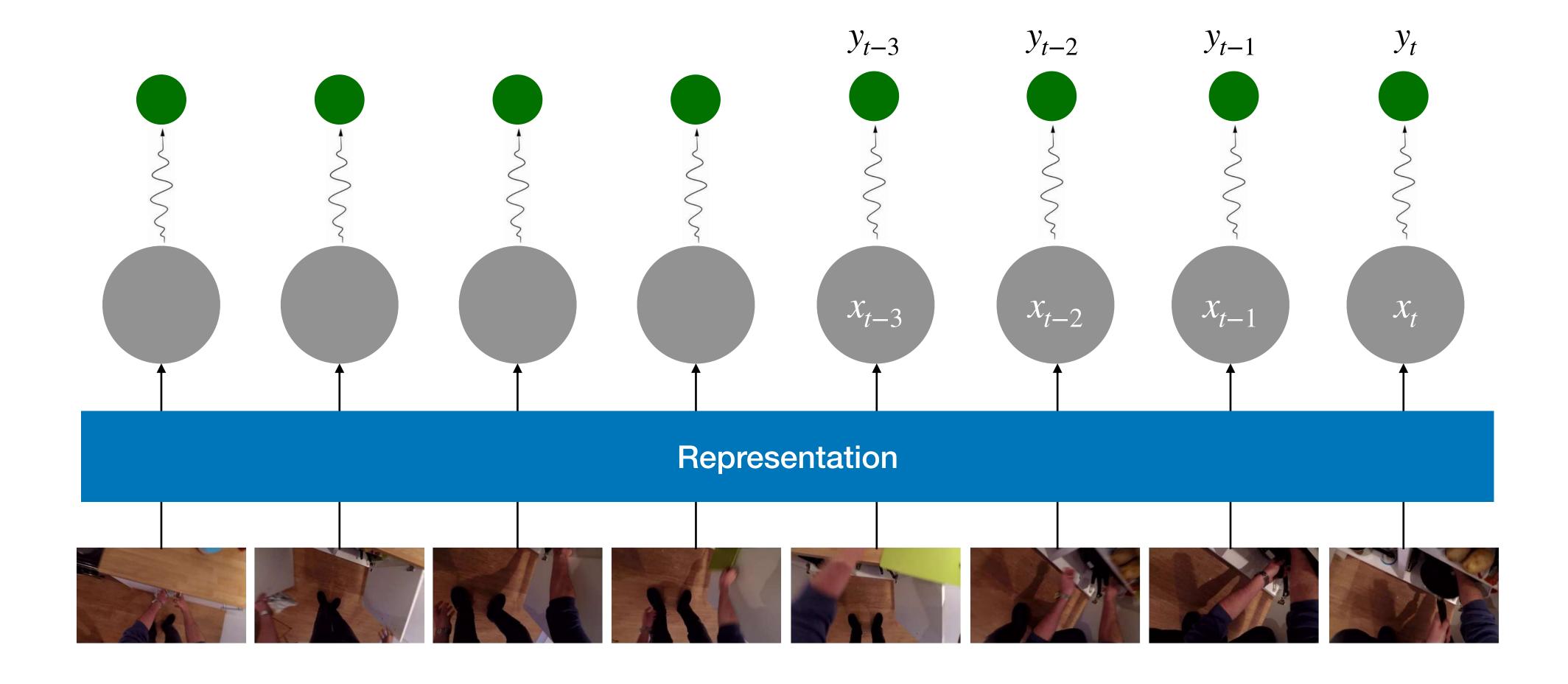


Consider a differential equations:

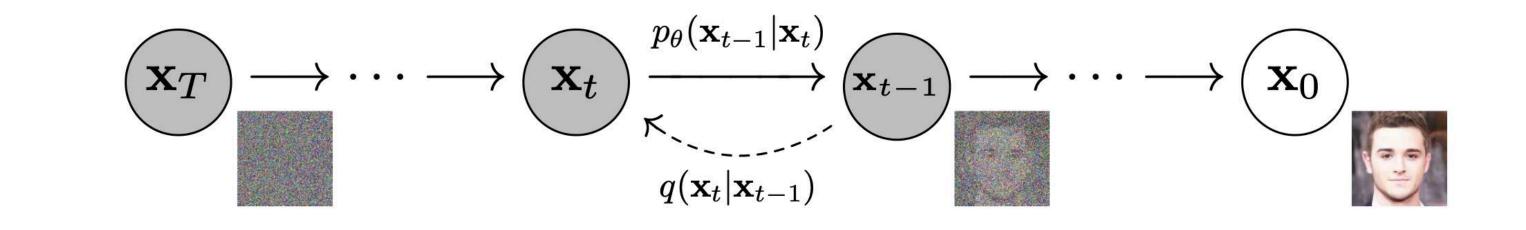
$$\frac{\partial x}{\partial t} = f(x)$$

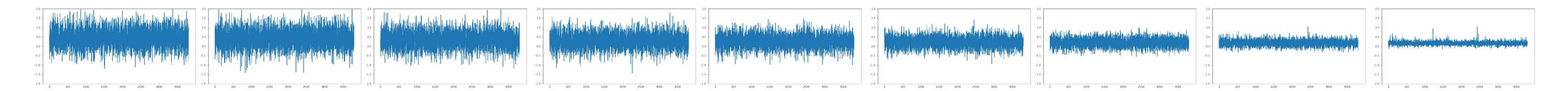
•
$$x(t_0) = x_0$$

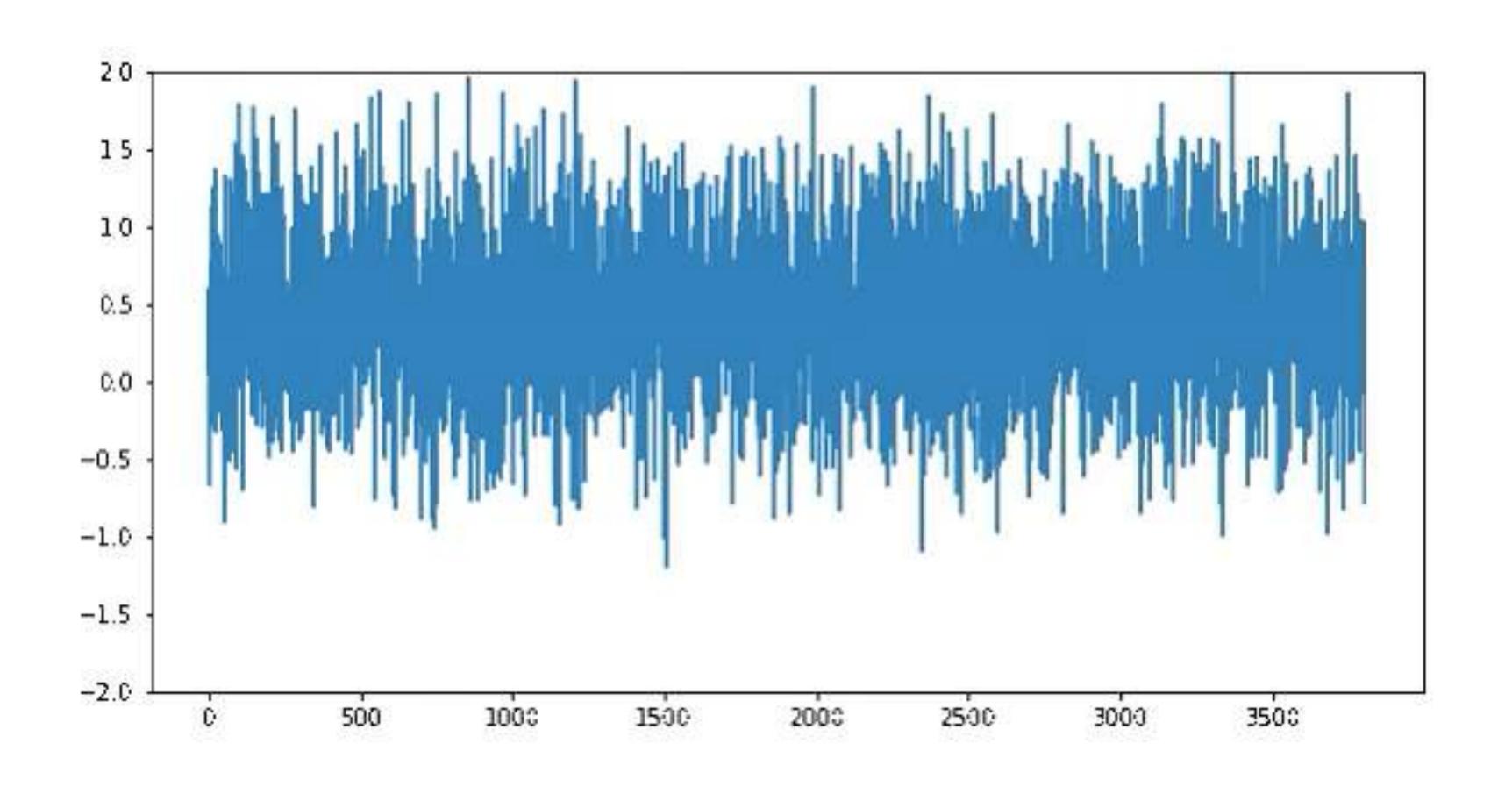
- The Picard-Lindelof theorem states that at least locally, if f is Lipschitz, there exists a unique solution x(t). Which implies:
 - $x(t + \partial t) = x(t) + \partial t \cdot f(x(t))$.



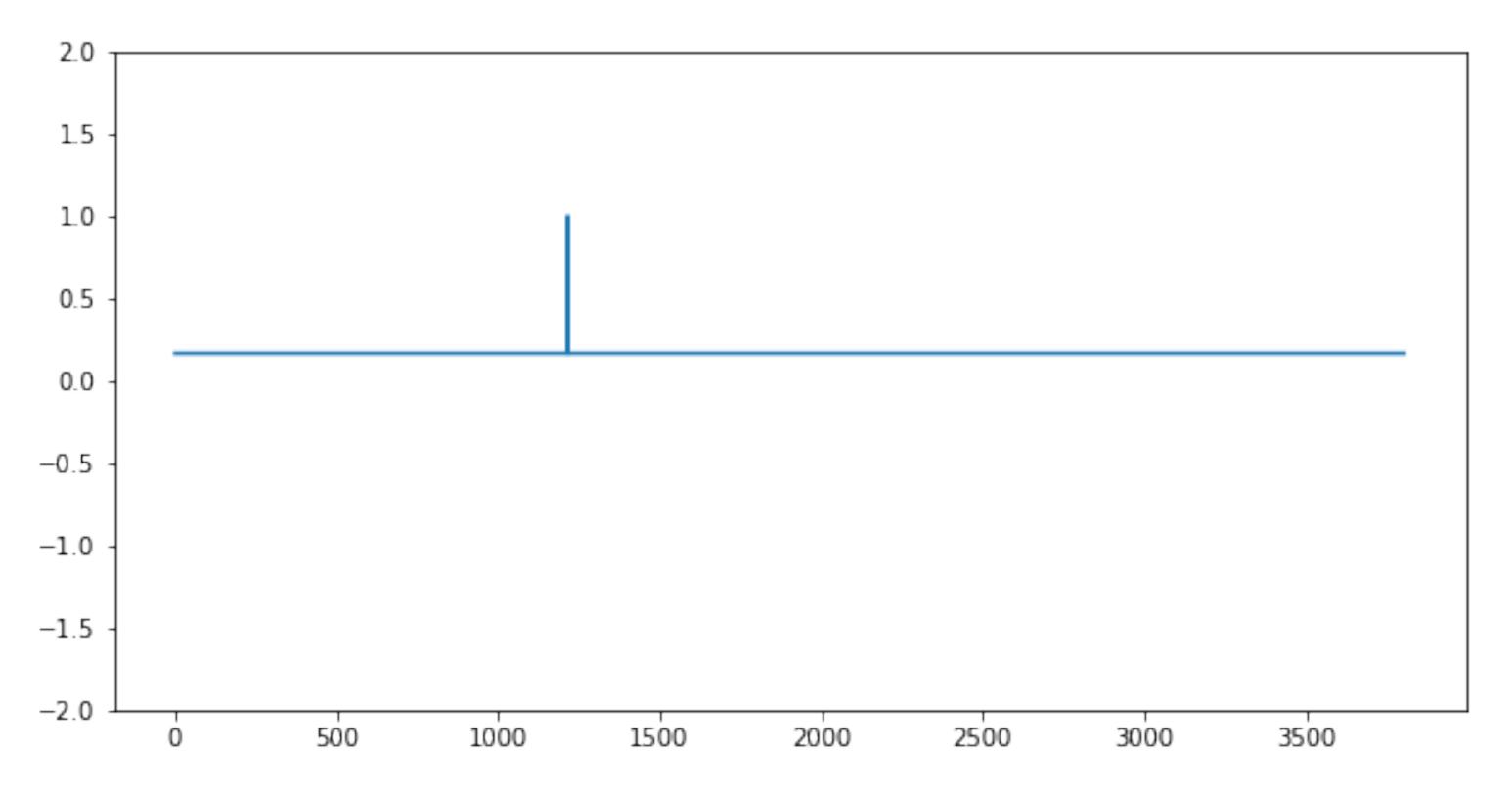
Integrating Video Action Anticipation with Diffusion Models







Prediction

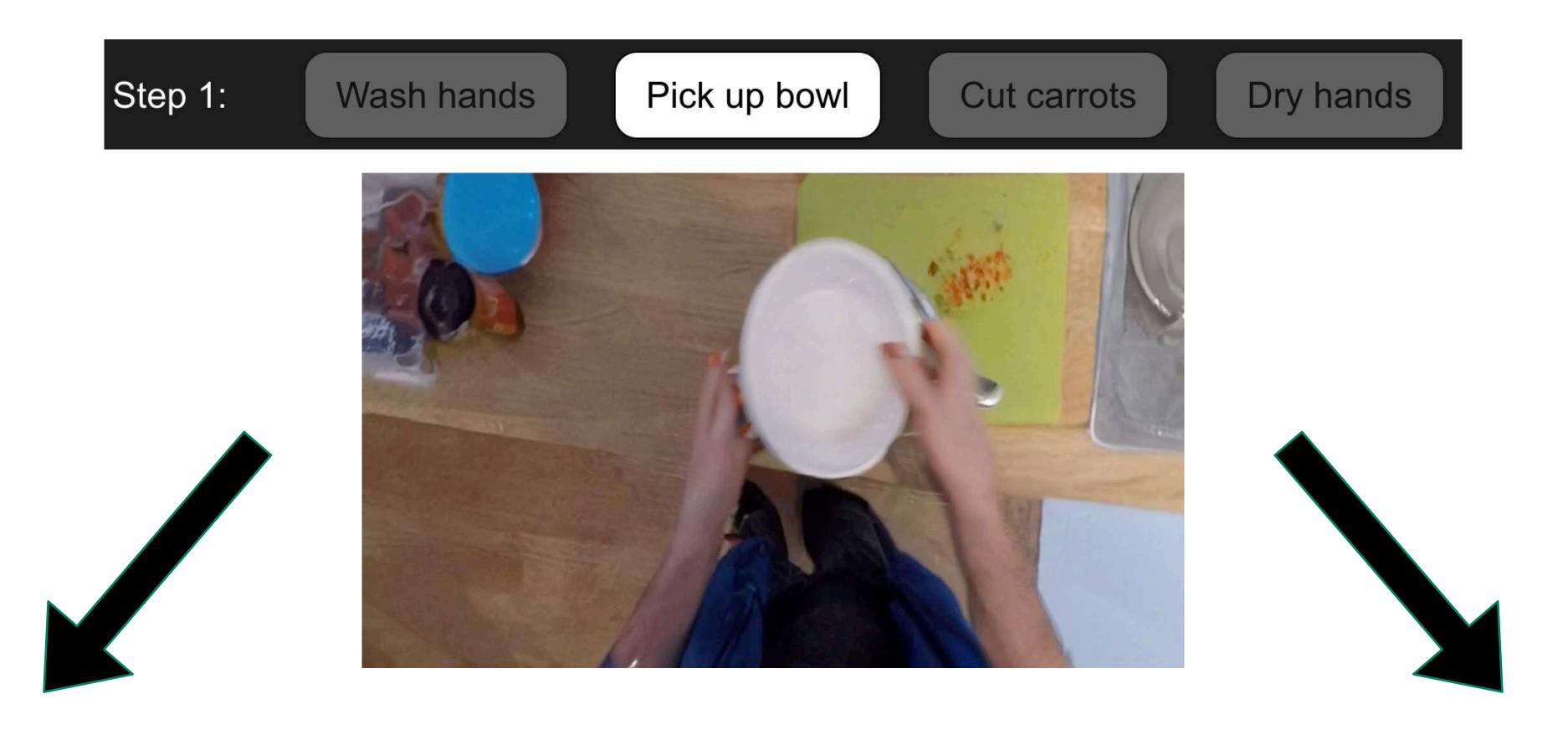


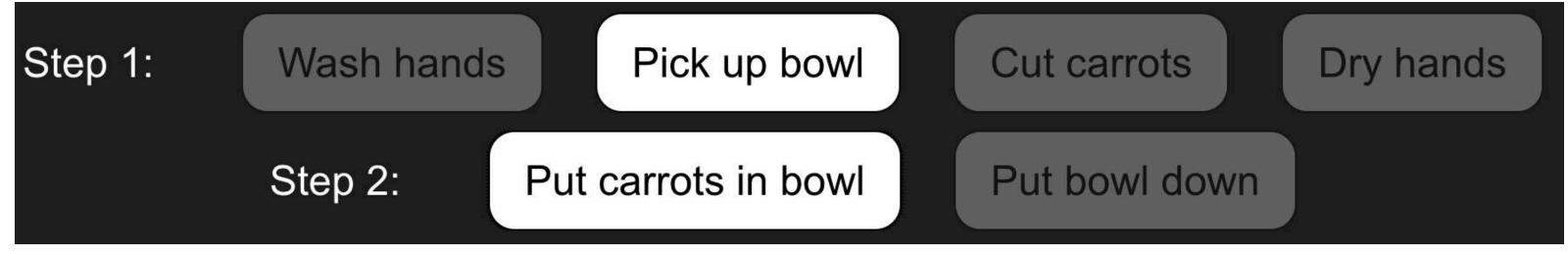
Ground Truth

Manipulate the Environments (GAIA-1, Wayve)

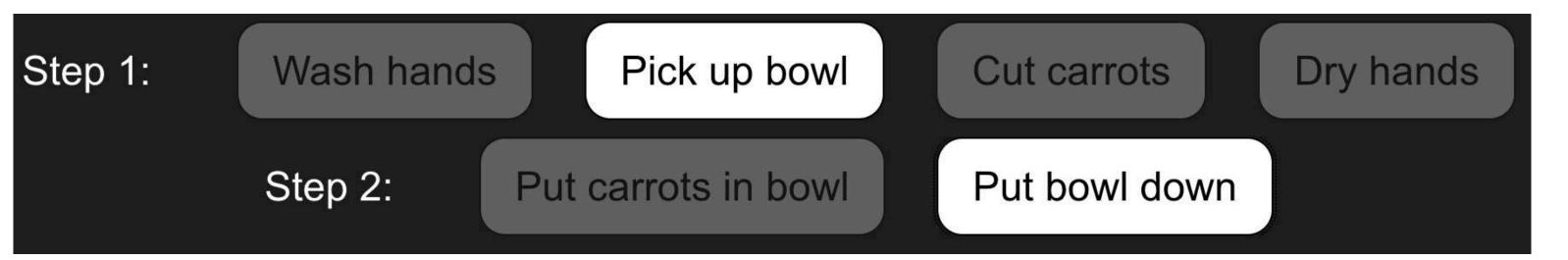


Analysis the Video Action Causality (UniSim, Yang et al. 2023)













Introduction

Video Action Anticipation

Introduction to Diffusion Models

Synthesis and Future Directions

Conclusion and Discussion

Conclusion and Discussion

- We've discussed the evolution of generative models, noting the progression from large language models (LLMs) to the current research focus on vision-language models (VLMs).
- The next frontier in the vision domain is video, which presents a significant challenge due to its inherent spatial and temporal complexities.
- The diffusion model operates on a discrete interpretation of Langevin Dynamics.
- The diffusion model is capable of synthesizing high-fidelity modalities, potentially serving the substantial data requirements of foundational models.

Can Generative Al and Video Analysis Redefine the Future Vision?

