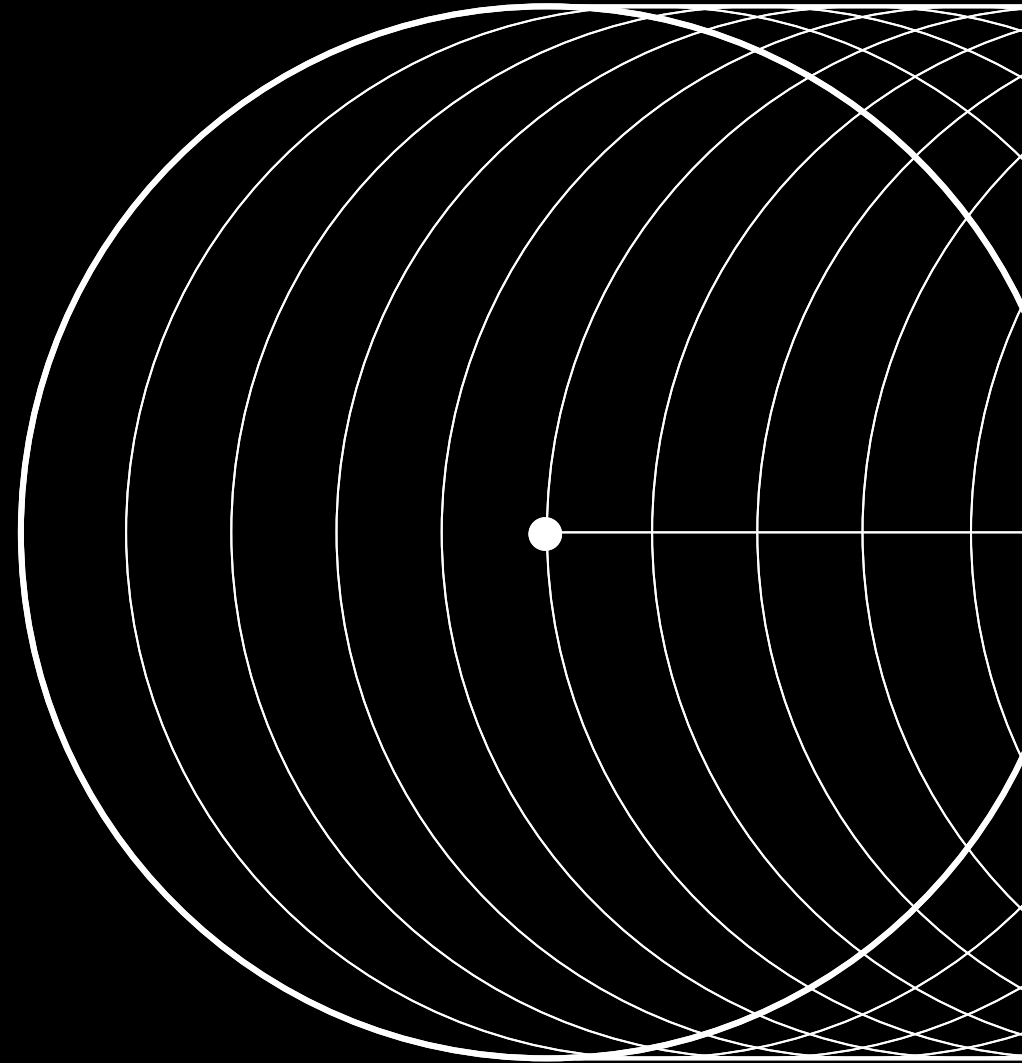# From Noise to Narrative: The Evolution of Visual Generative Models

**Sergey Kastryulin**

Research Scientist

# Agenda

# Visual Generative Modeling



**Unconditional**

# Visual Generative Modeling



Unconditional

Class-conditional

# Visual Generative Modeling



"Cyberpunk girl"

**Unconditional**            **Class-conditional**            **Text-conditional**

YandexART

# YaART: Yet Another ART Rendering Technology
## KDD'25 and Beyond

YandexART

ALCHEMIST
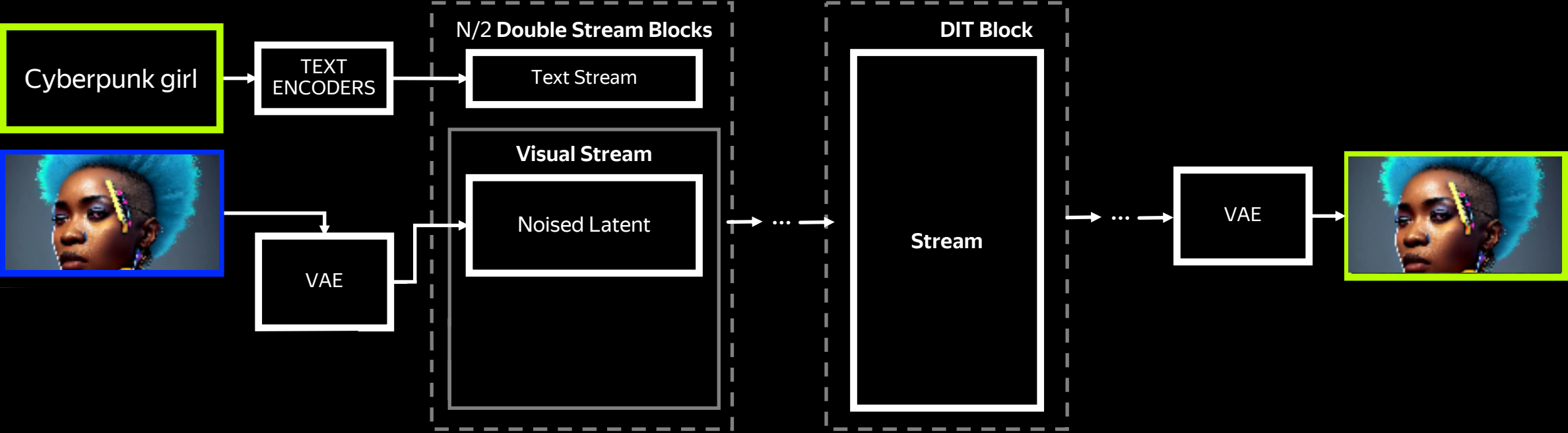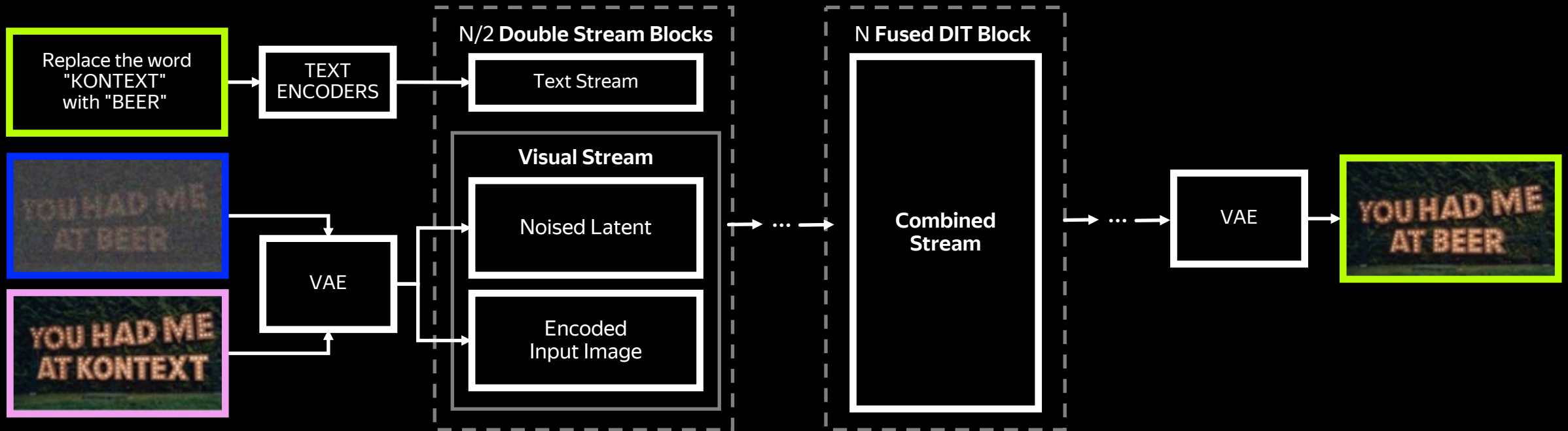
NeurIPS'25
Fall into ML — 16:00, poster #39)

7

# Conditioning on text

In general

# Conditioning on text and image

**Flux.1 Kontext**



S. Batifol et. al., "FLUX.1 Kontext: Flow Matching for In-Context Image Generation and Editing in Latent Space", 2025
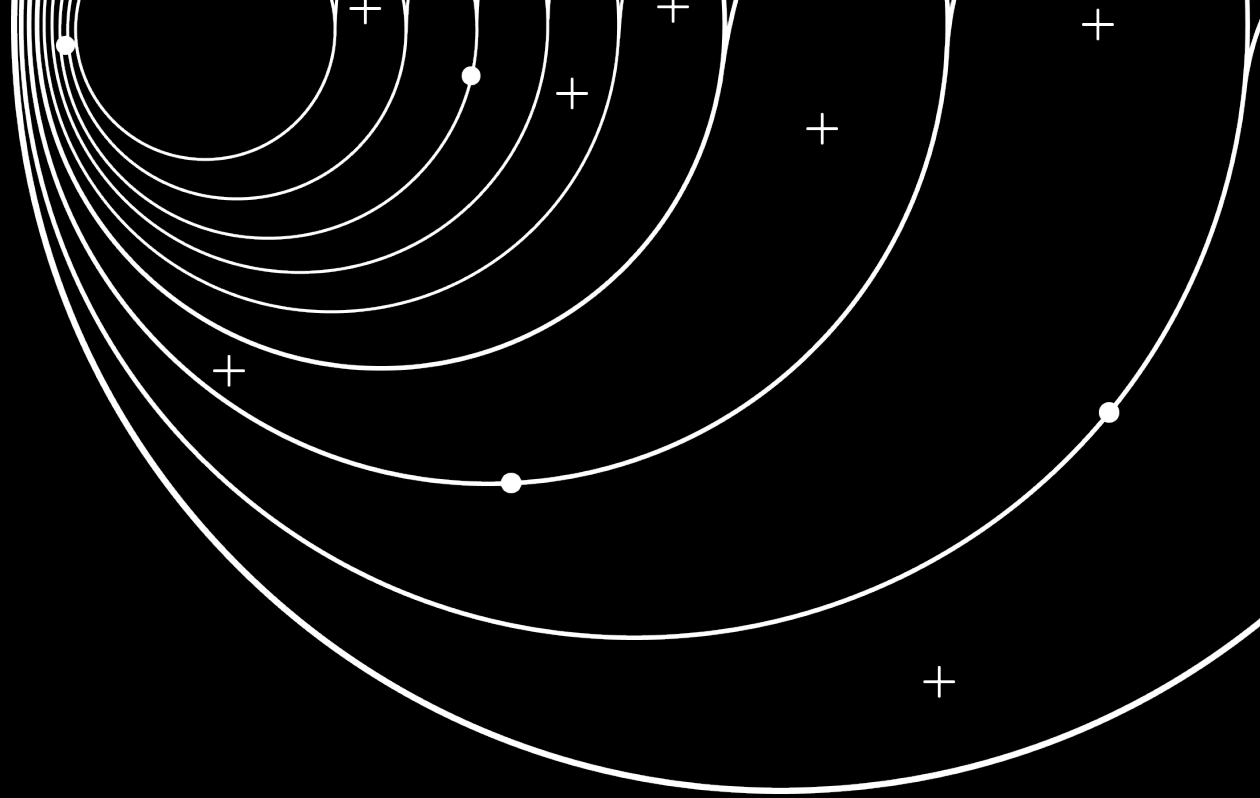
# Conditioning on text and image



Restore photo

Add sports costume

# In reality, we want a dialog

# In reality we want a dialog

Why?

**01** Generate new content

**02** Edit existing content

**03** Ask questions

**04** Work with several images

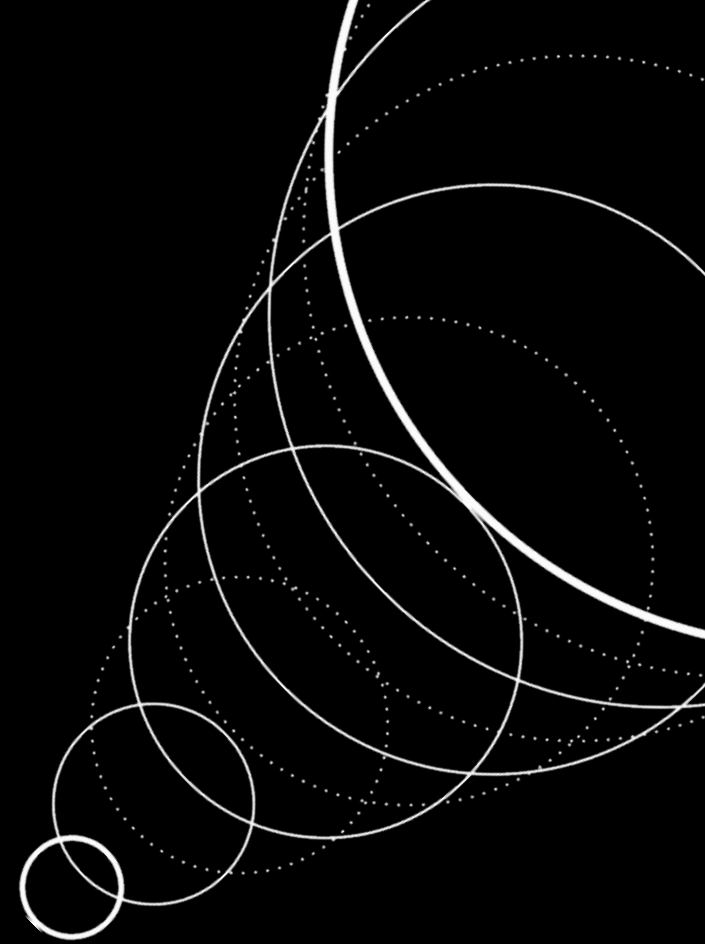**05** Support context

**06** Bonuses from different modalities*

---

* J. Zhang et. al., "Are Unified Vision-Language Models Necessary: Generalization Across Understanding and Generation", 2025

# Dialog means unification

## Unify two main generative worlds

- Inherently continuous (visual): images, video, 3D
- Inherently discrete (textual): text, code, math
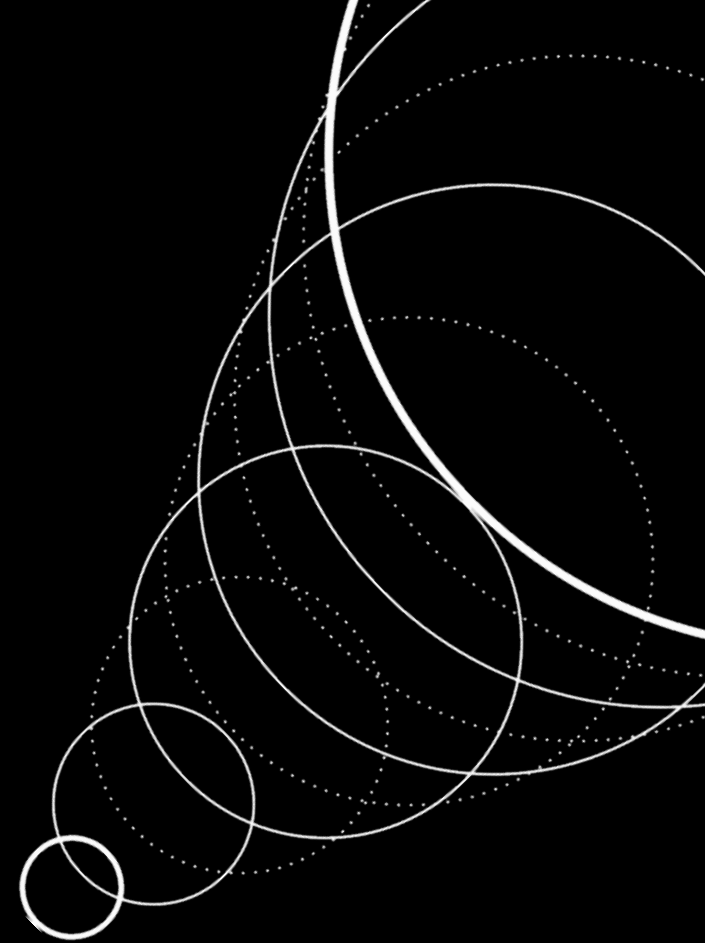
# Dialog means unification
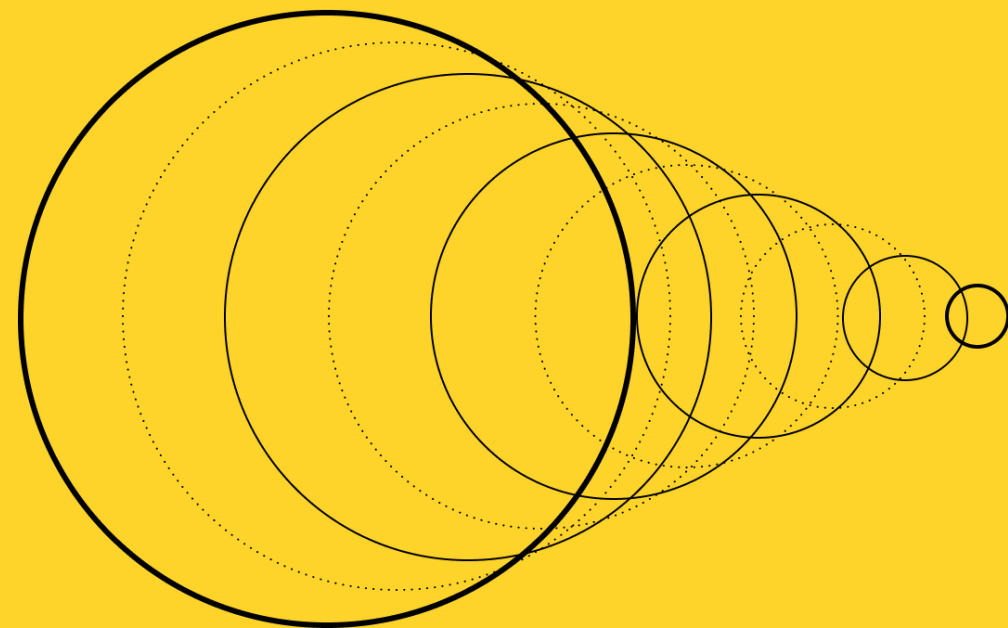
## Unify two main generative worlds

- Inherently continuous (visual): images, video, 3D
- Inherently discrete (textual): text, code, math

## Currently we can model

- text-to-text — LLM
- text-to-image — all we discussed above
- image-to-text — VLM
- text+image-to-image — Editing

Yandex Research
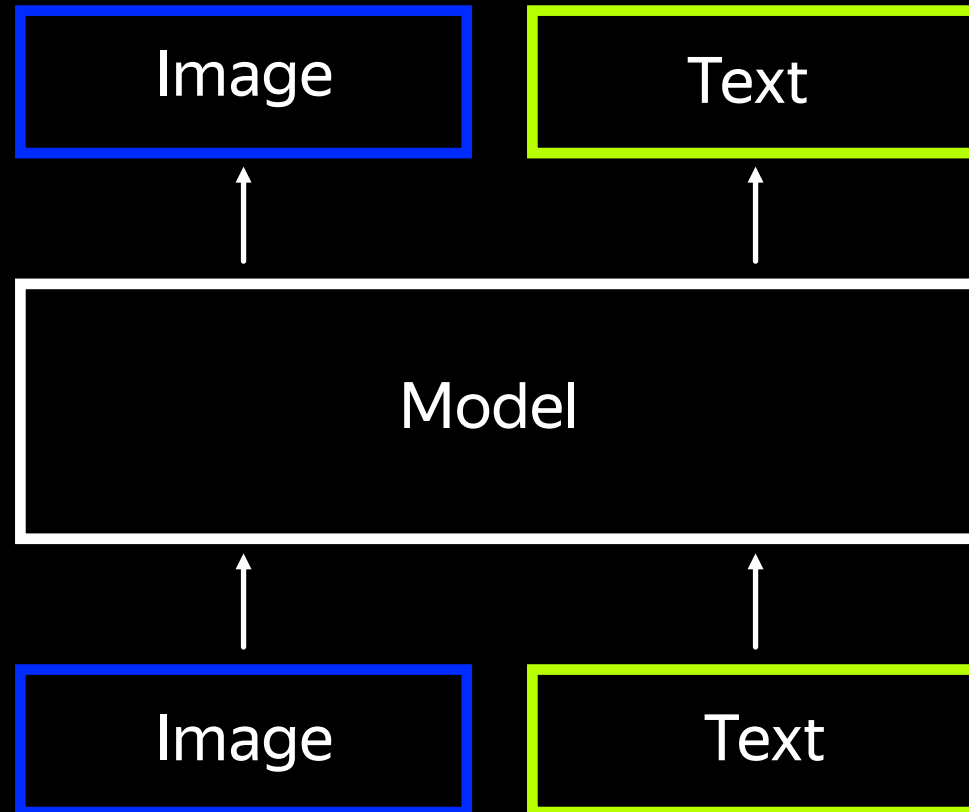
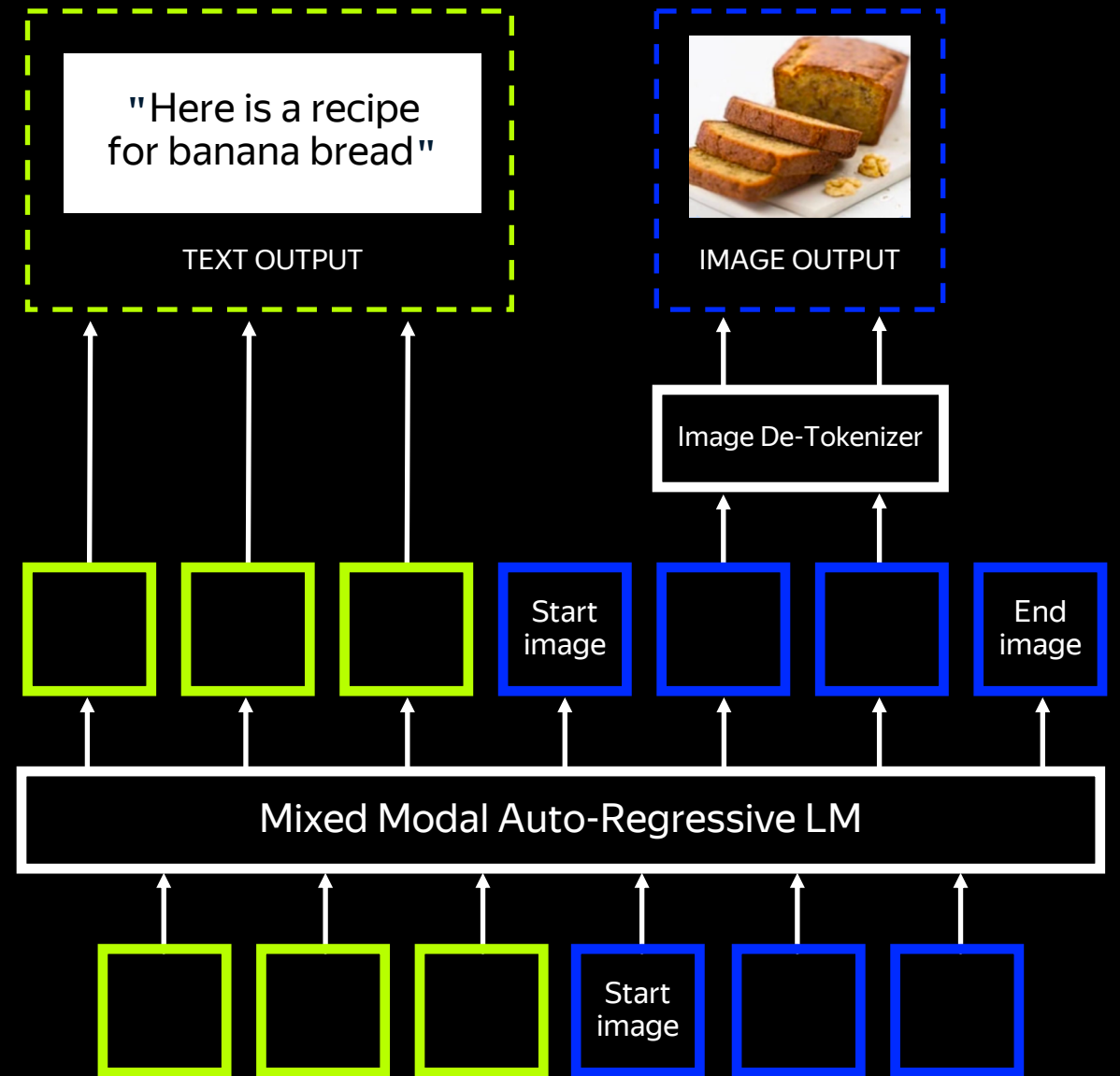# Towards Unified Generative Models

# What kind of model would help us?

# Naive approach

"Here is a recipe for banana bread"

TEXT OUTPUT

IMAGE OUTPUT

Image De-Tokenizer

Start image

End image

Mixed Modal Auto-Regressive LM

Start image

# Naive approach

$\oplus$ Uniform implementation

$\ominus$ Raster order is not native to images

$\ominus$ Discretisation ruins image quality



"Here is a recipe for banana bread"

TEXT OUTPUT

IMAGE OUTPUT

Image De-Tokenizer

Start image

End image

Mixed Modal Auto-Regressive LM

Start image

# Naive approach

$\oplus$ Uniform implementation

$\ominus$ Next-token prediction is not native to images

$\ominus$ Discretisation ruins image quality

sic!

"Here is a recipe for banana bread"

TEXT OUTPUT

IMAGE OUTPUT

Image De-Tokenizer

Start image

End image

Mixed Modal Auto-Regressive LM

Start image
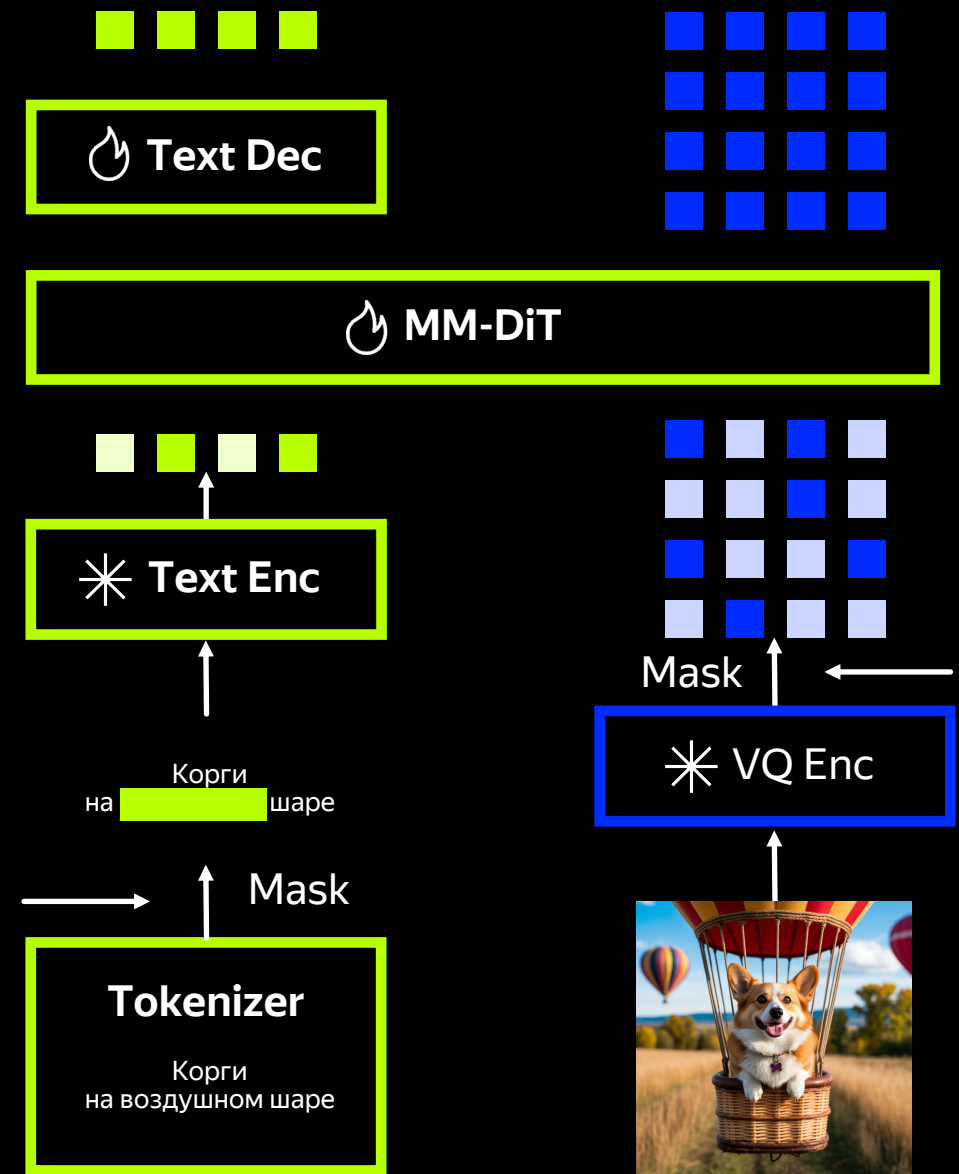
Start image

A. Ramesh et. al., "Zero-Shot Text-to-Image Generation" (aka DALL-E 1), 2021

19

# Muddit

(+) Uniform implementation

(−) Discrete diffusion for text data*

(−) Discrete tokens for image data

🔥 **Text Dec**

🔥 **MM-DiT**

✳ **Text Enc**

Корги
на ▮▮▮▮ шаре

→ Mask

**Tokenizer**

Корги
на воздушном шаре

Mask ↑ ←

✳ VQ Enc

Q. Shi et. al., "Muddit: Liberating Generation Beyond Text-to-Image with a Unified Discrete Diffusion Model", 2025

# Transfusion

⊕ Unified representation in a single transformer

⊕ Causal attn for text, bidirectional for images

⊕ Continuous encoding of images

VAE Decoder

Linear ⬜ or 🟪 U-Net Up

Transformer

Linear ⬜ or 🟪 U-Net Down

Noising

VAE Encoder

$$\mathcal{L}_{Transfusion} = \mathcal{L}_{LM} + \lambda \times \mathcal{L}_{DDPM}$$

C. Zhou et. al., "Transfusion: Predict the Next Token and Diffuse Images with One Multi-Modal Model", 2024

# MAR, (Uni)Fluid

(+) Use pre-trained image encoders (VAEs)

(+) Continuous representation of images

(−) Long sampling

**Generation**

**Per Token Diffusion Head**

**Understanding**

**Per Token Classification Head**

**Unified AR Transformer**

<BOI> <BOS>

Text Tokenizer

Image Tokenizer (VAE)

Image Encoder (SigLIP)

Text Tokenizer

Text

Image

Image

Text

T. Li et. al., "Autoregressive Image Generation without Vector Quantization", 2024

# NexusGen | Qwen Image



+ Re-use pre-trained diffusion denoiser

+ Should be cheaper to train, right?*

− Hard to align with good quality

**UnPatchify**

**MMDiT Block**

×N

**MMDiT Block**

✳ **Qwen2.5 VL**

System promt    User promt

Three tall coconut trees, two of which are located in the center of the picture and one on the right edge. Two of the coconut trees are covered in golden coconut fruits under a clear blue sky, presenting a peaceful and tropical natural scene.

Patchify

t

Noize

✳ VAE Encoder

C. Wu et. al., "Qwen-Image Technical Report", 2025

C. Wu et. al., "Qwen-Image Technical Report", 2025

# How hard is it to train the connector?

*"a pair of cherries, dressed in a delicate ballerina outfit"*





- Images generated with **NexusGen** have low quality
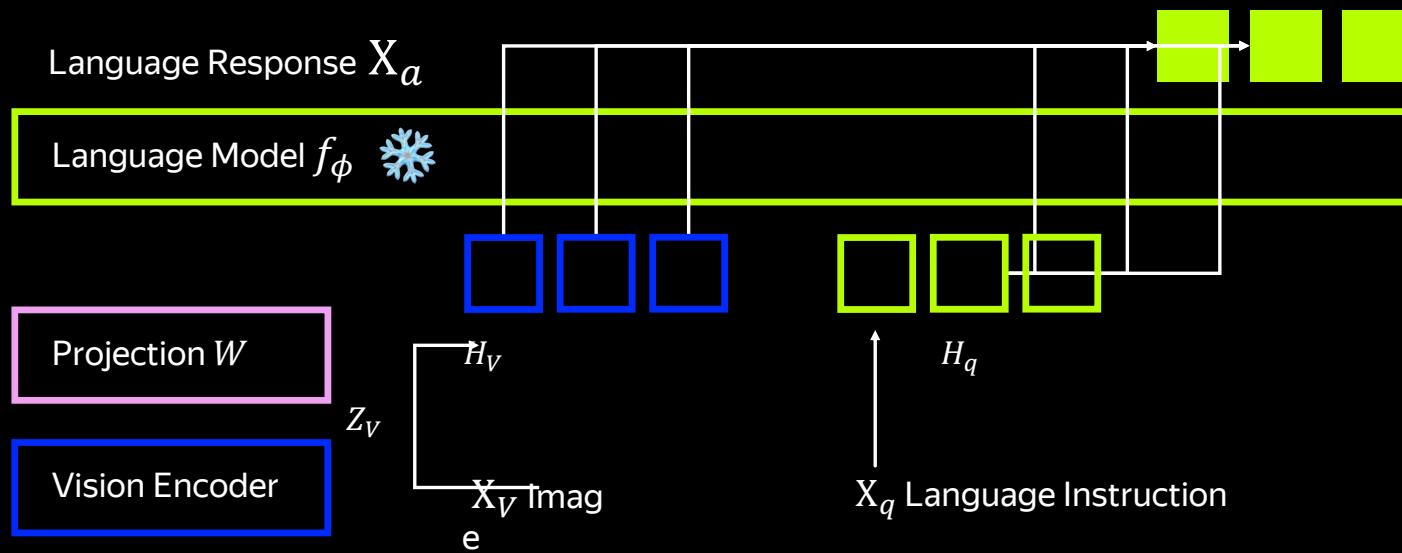
- Did not train from scratch

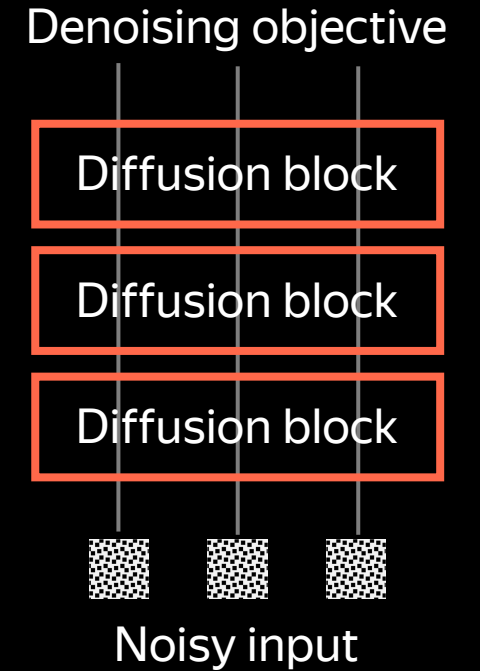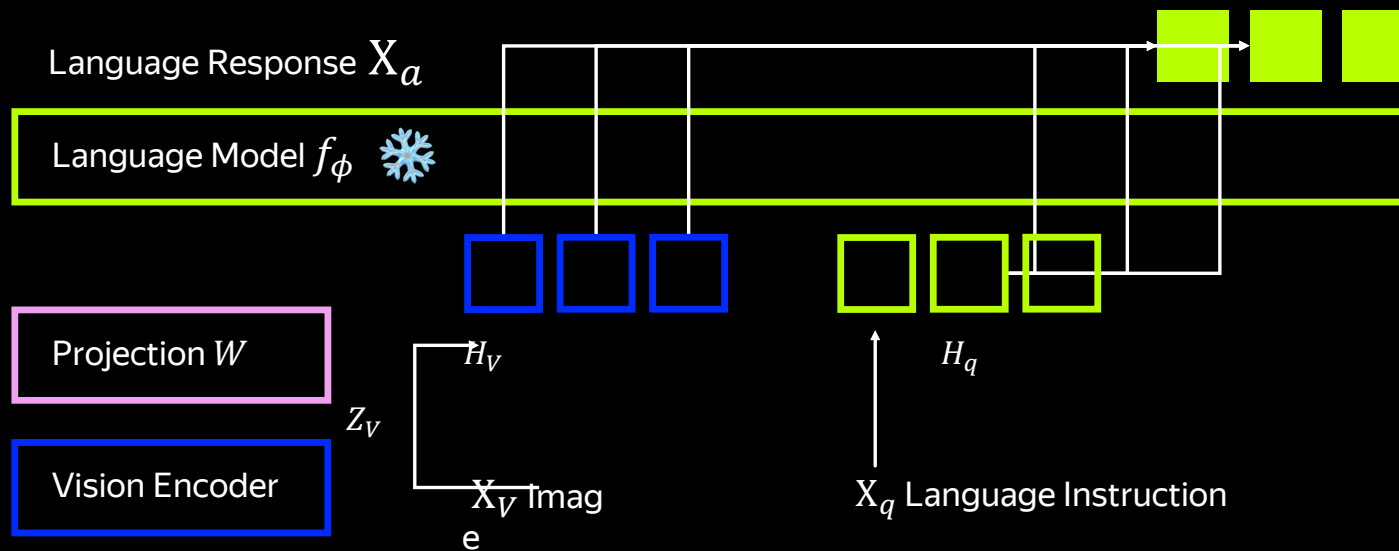- Images generated with **Qwen Image** are much better

- Trained from scratch

* Very much a hand-waving argumentation, lots of other variables are not aligned between setups
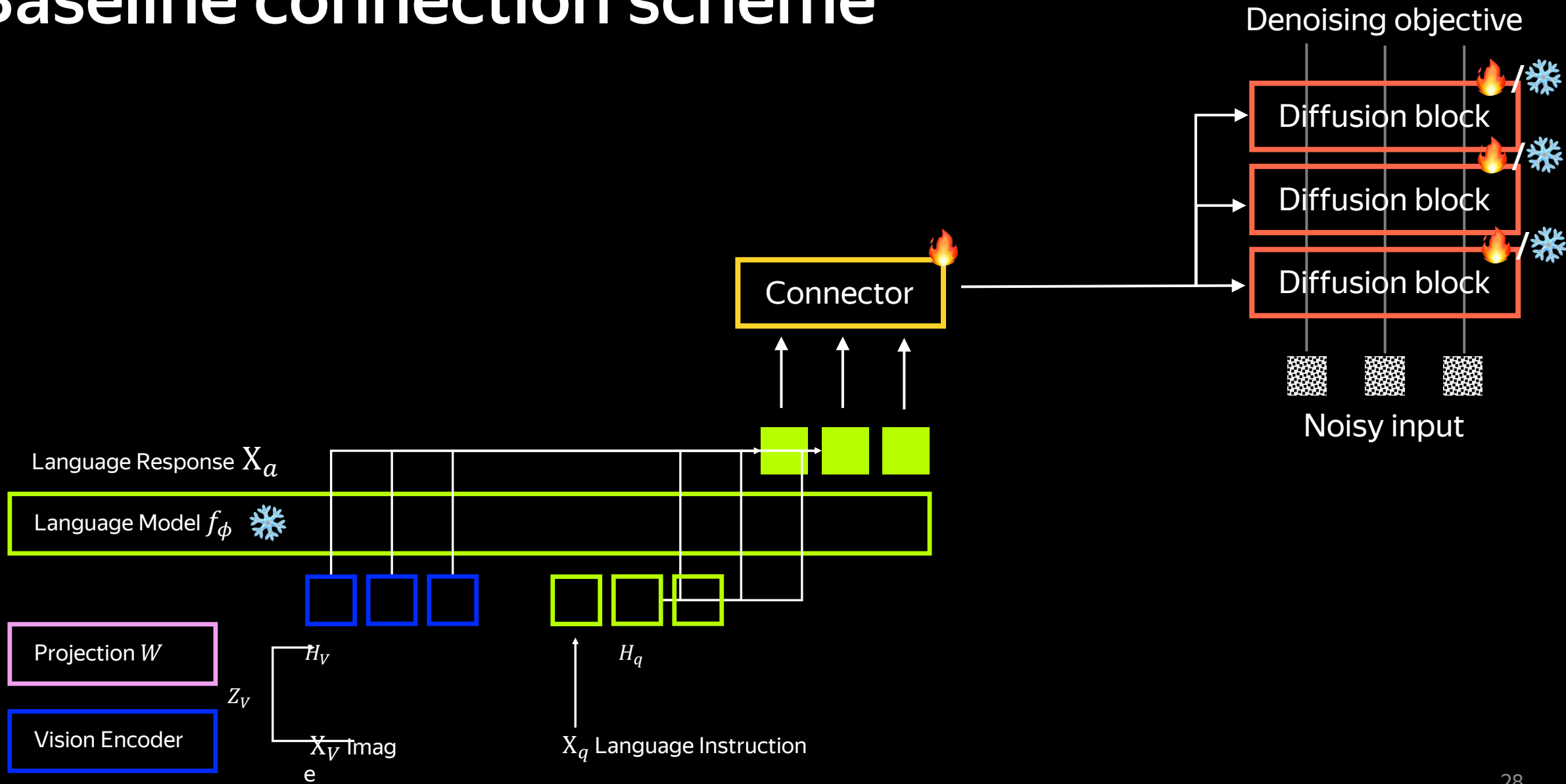
# Baseline connection scheme

Language Response $\mathrm{X}_a$
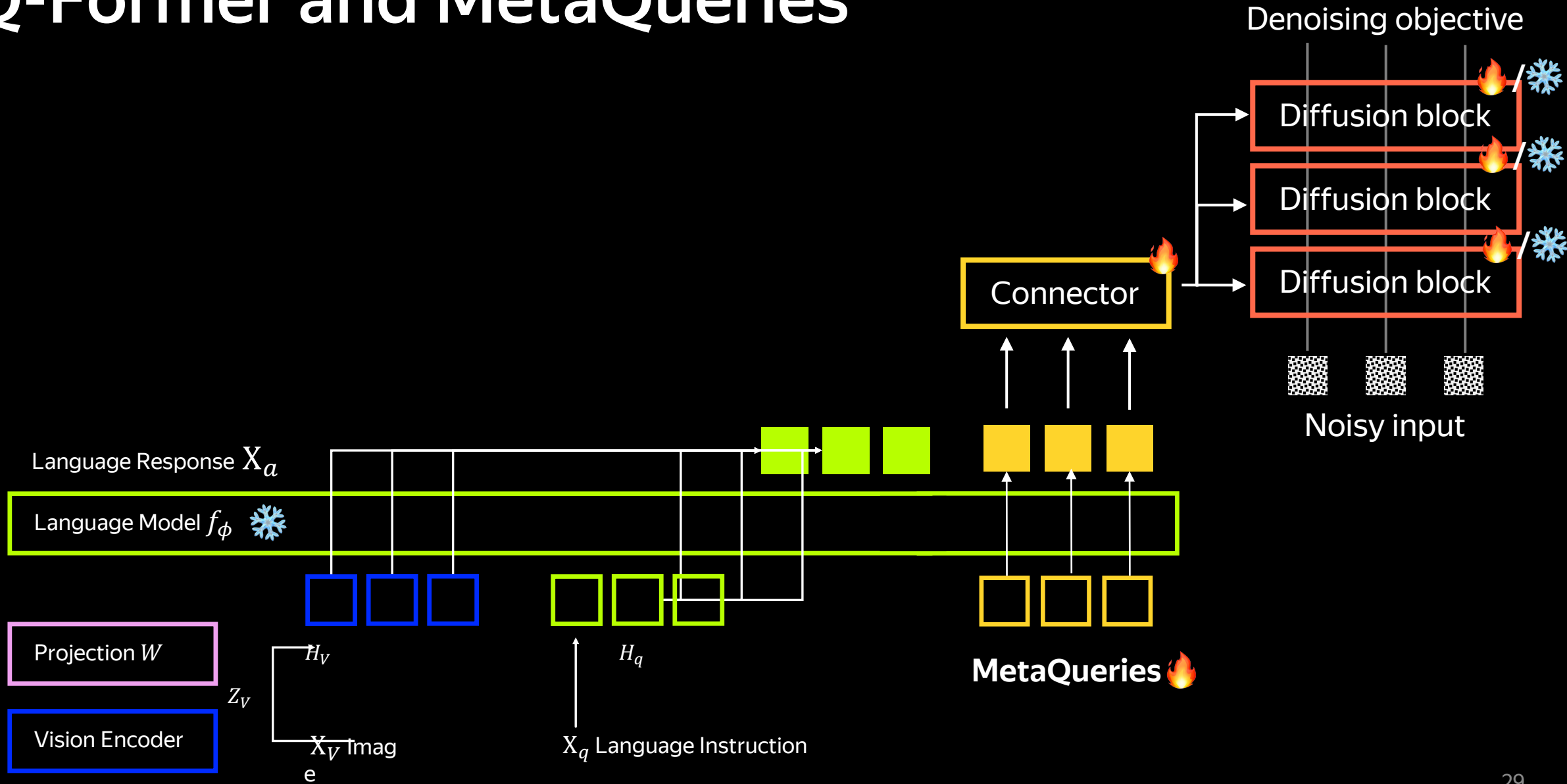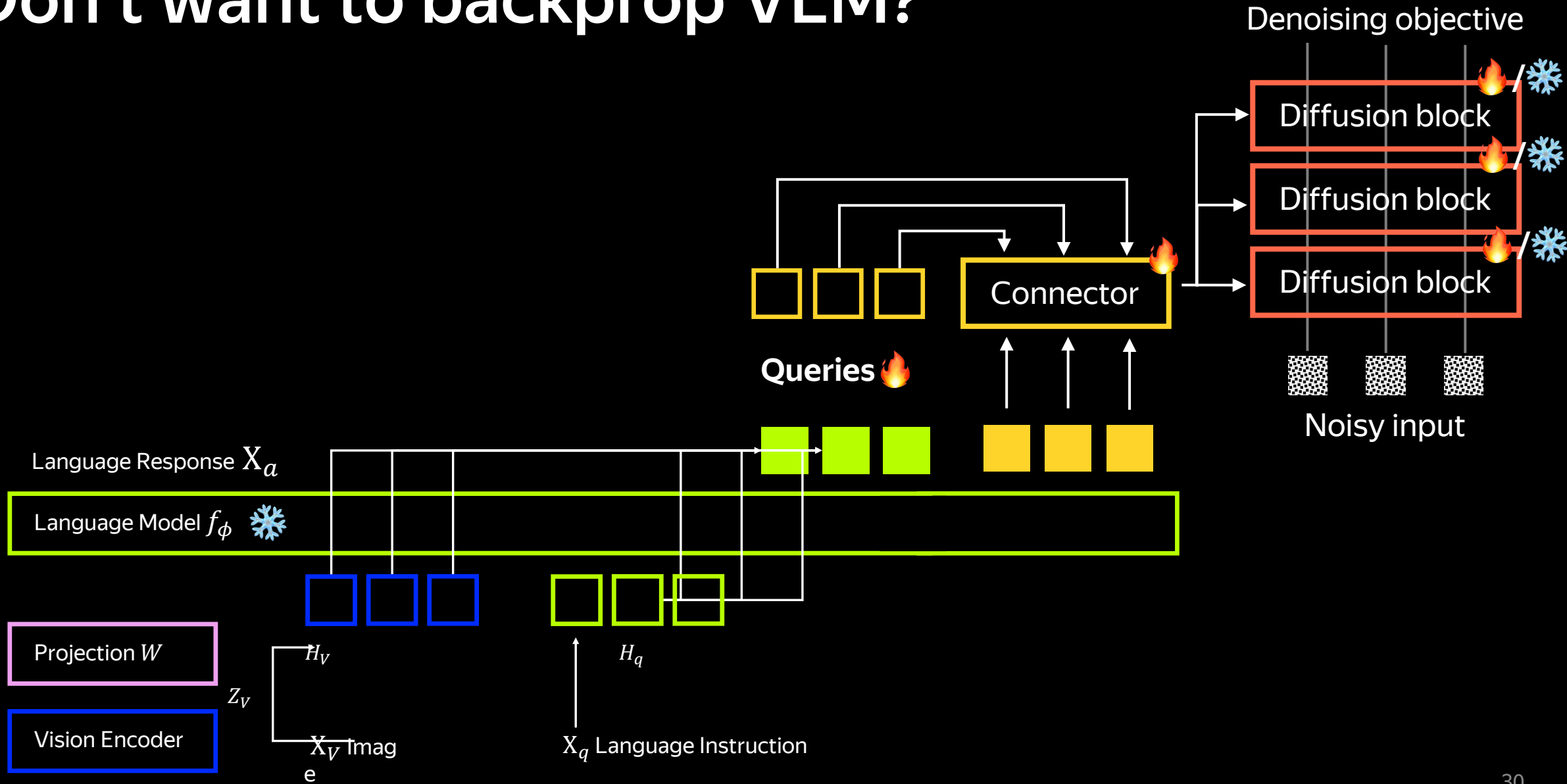
Language Model $f_\phi$ ❄️

Projection $W$

$Z_V$

Vision Encoder

$H_V$

$H_q$

$\mathrm{X}_V$ Image

$\mathrm{X}_q$ Language Instruction

# Baseline connection scheme

Denoising objective

Diffusion block

Diffusion block

Diffusion block

Noisy input

Language Response $\mathrm{X}_a$

Language Model $f_\phi$

Projection $W$

Vision Encoder

$Z_V$

$H_V$

$\mathrm{X}_V$ Image

$H_q$

$\mathrm{X}_q$ Language Instruction

# Baseline connection scheme

Denoising objective

Diffusion block 🔥/❄️

Diffusion block 🔥/❄️

Diffusion block 🔥/❄️

Connector 🔥

Noisy input

Language Response $X_a$

Language Model $f_\phi$ ❄️

Projection $W$

$H_V$

$Z_V$

Vision Encoder

$X_V$ Image

$H_q$

$X_q$ Language Instruction

# Q-Former and MetaQueries



Denoising objective

Diffusion block 🔥/❄️

Diffusion block 🔥/❄️

Diffusion block 🔥/❄️

Connector 🔥

Noisy input

Language Response $X_a$

Language Model $f_\phi$ ❄️

Projection $W$

Vision Encoder

$H_V$

$Z_V$

$X_V$ Image

$H_q$

$X_q$ Language Instruction

MetaQueries 🔥

X. Pan et. al., "Transfer between Modalities with MetaQueries", 2025

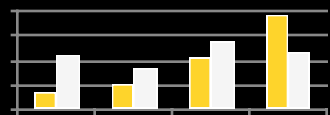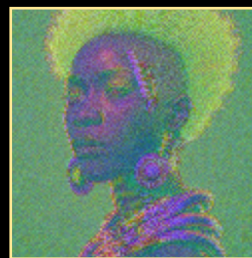# Don't want to backprop VLM?

# Recap



Continuous

Discrete

Images representation

| | |
|---|---|
| **Image** | Text |

Model

| | |
|---|---|
| **Image** | Text |

# Recap



Diffusion     Autoregression

Loss

Image     Text

Model

Image     Text

Continuous

Images representation

Discrete
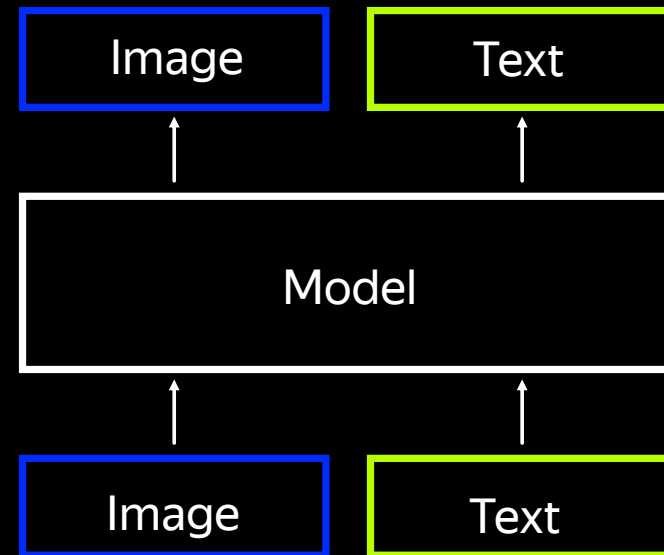
# Recap

**Image gen module**



**External**                                    **Integrated**

# Recap

**Image gen module**



Image

Denoiser        Text

Connector? →

Model

Image        Text

**External**

Image        Text

Model

Image        Text

**Integrated**

# BAGEL

Language response

FFN

Self-attention

QKV

Text tokenizer          Und encoder

C. Deng et. al., "Emerging Properties in Unified Multimodal Pretraining", 2025

# BAGEL

Language response      Language response

C. Deng et. al., "Emerging Properties in Unified Multimodal Pretraining", 2025

# BAGEL

C. Deng et. al., "Emerging Properties in Unified Multimodal Pretraining", 2025

# BAGEL

Language response      Image/Multi-image/Video

FFN      FFN

Multi-modal self-attention

QKV      QKV

Text tokenizer     Und encoder     Gen encoder     🤔 Integrated?

C. Deng et. al., "Emerging Properties in Unified Multimodal Pretraining", 2025
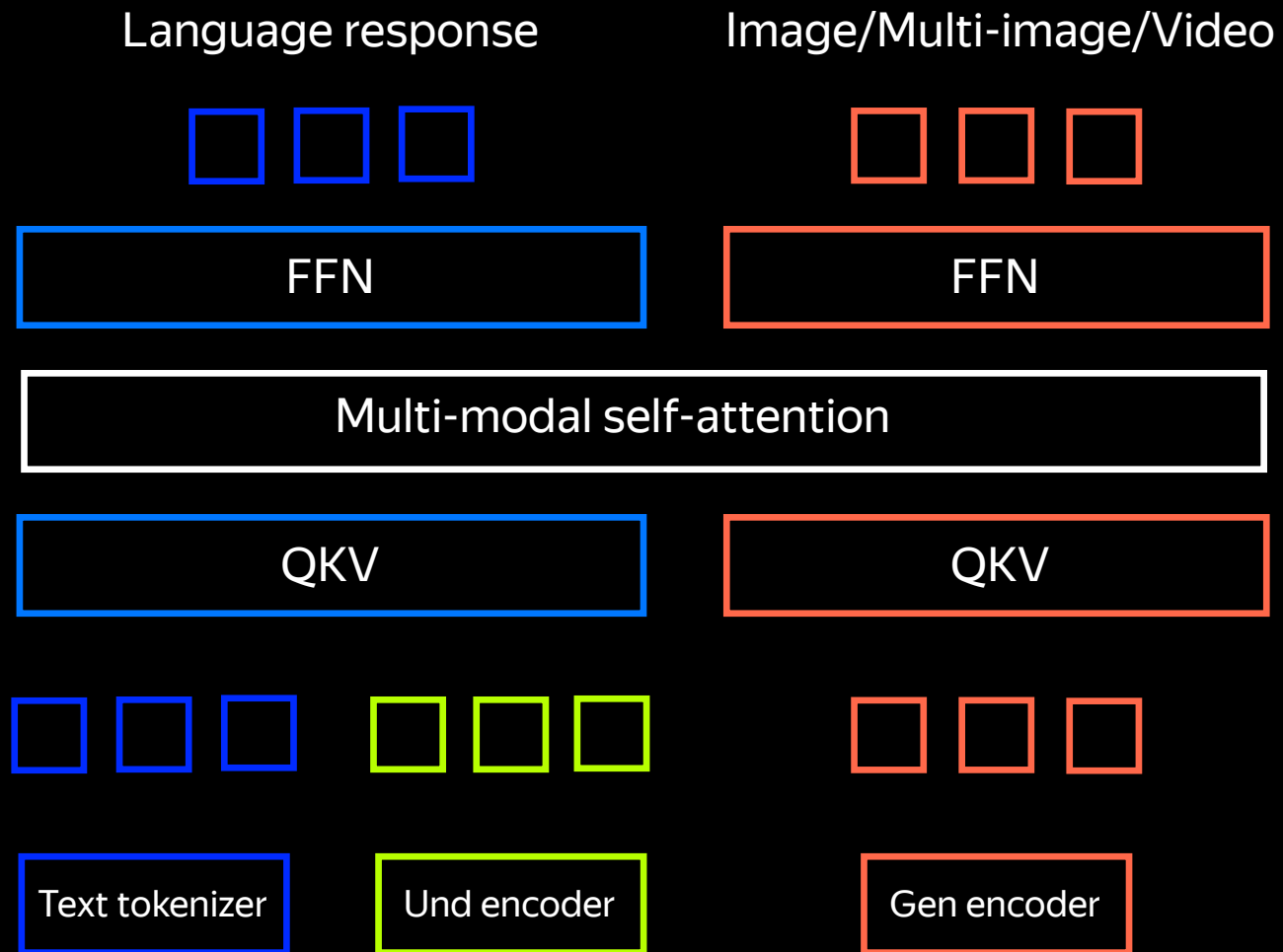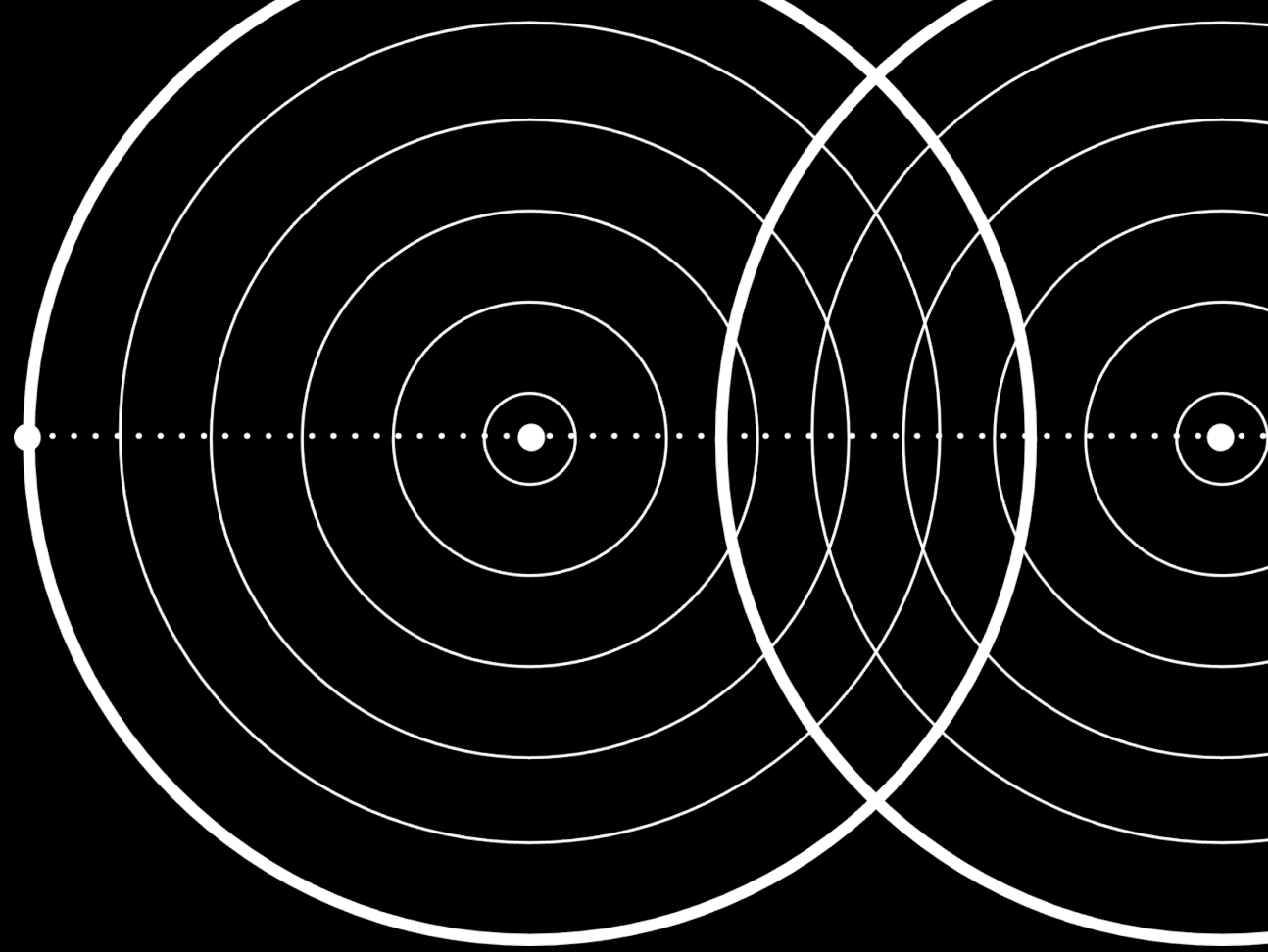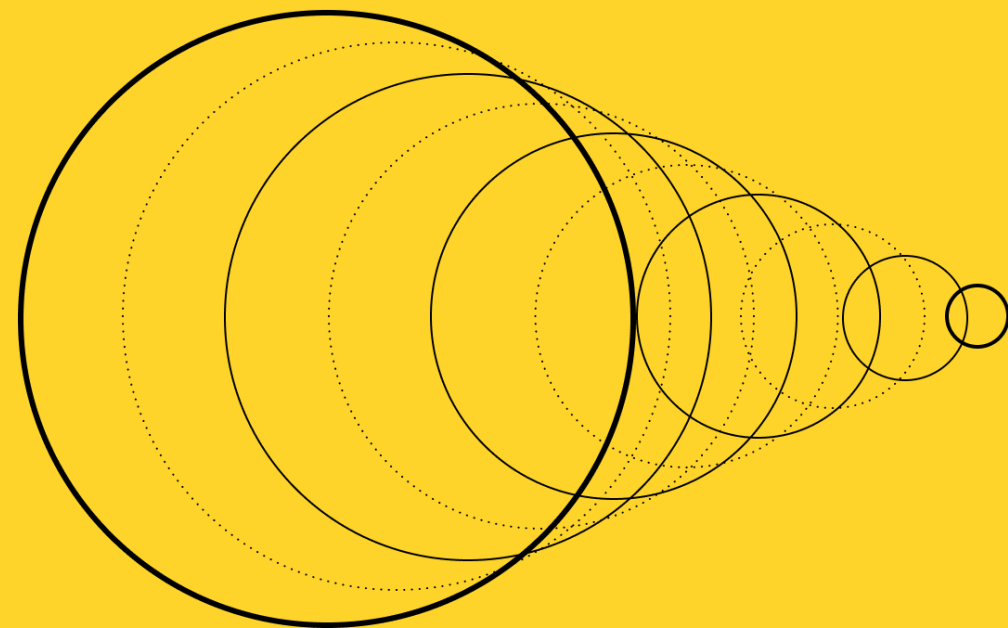
# Final recap

- Plenty of options to choose
- Options have ups and downs
- There is no clear winner yet

- Many open questions

# (Open) Questions

# How to Train

**Currently, we can model**

- text-to-text — LLM
- text-to-image — text-cond diffusion
- image-to-text — VLM
- text+image-to-image — text+img-cond diffusion

# How to Train

## Currently we can model

- text-to-text — LLM
- text-to-image — text-cond diffusion
- image-to-text — VLM
- text+image-to-image — text+img-cond diffusion



The moment you start you be like: how the hell do I combine all that?
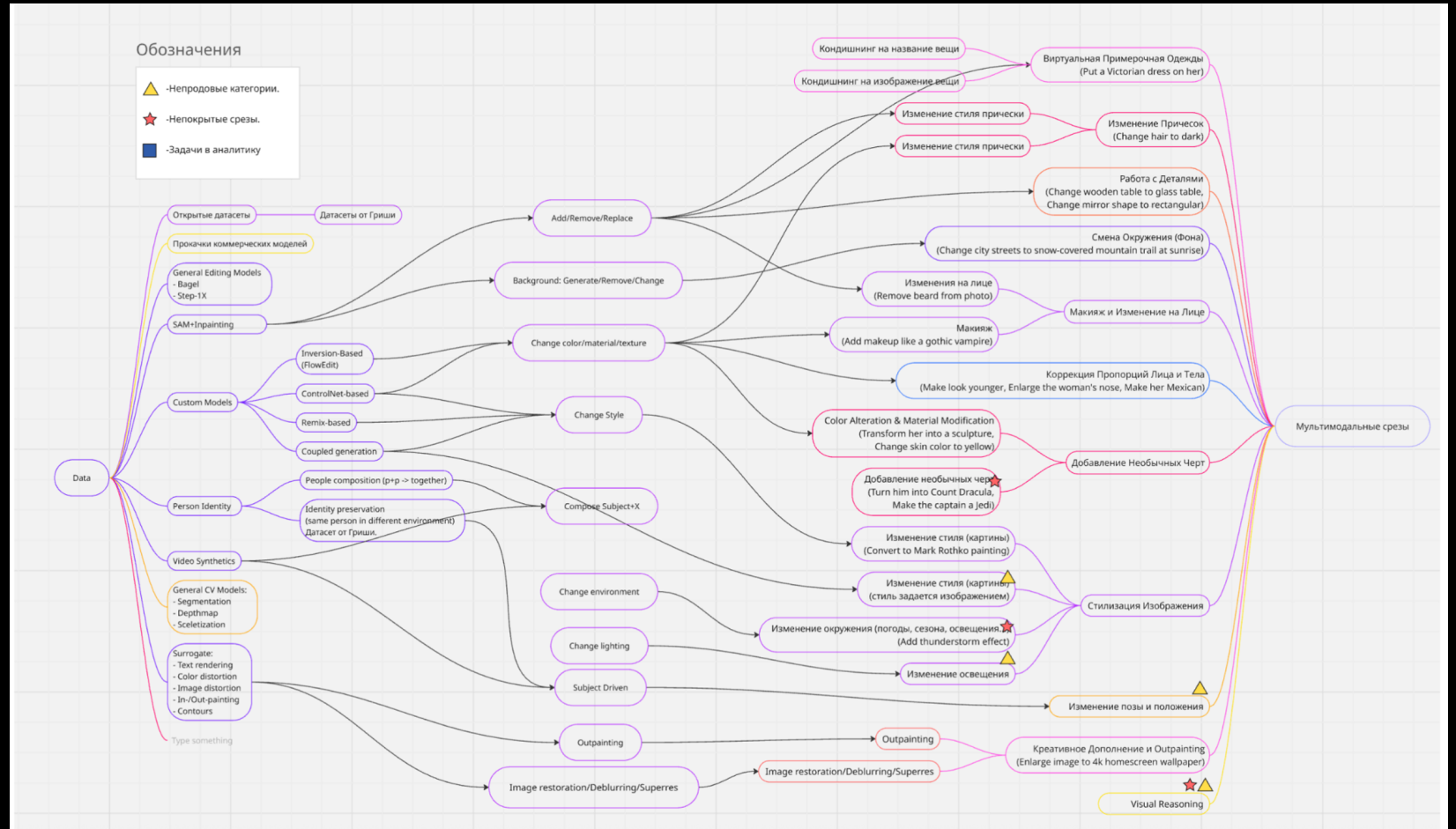
# Is there cross-modal knowledge transfer?

- Editing-only training collapses T2I

- Training on T2I + editing > just T2I or just editing

- Training on all 4 tasks is not worse than training on them individually

* all conclusions are from our internal experiments

# Adapted setup does magic

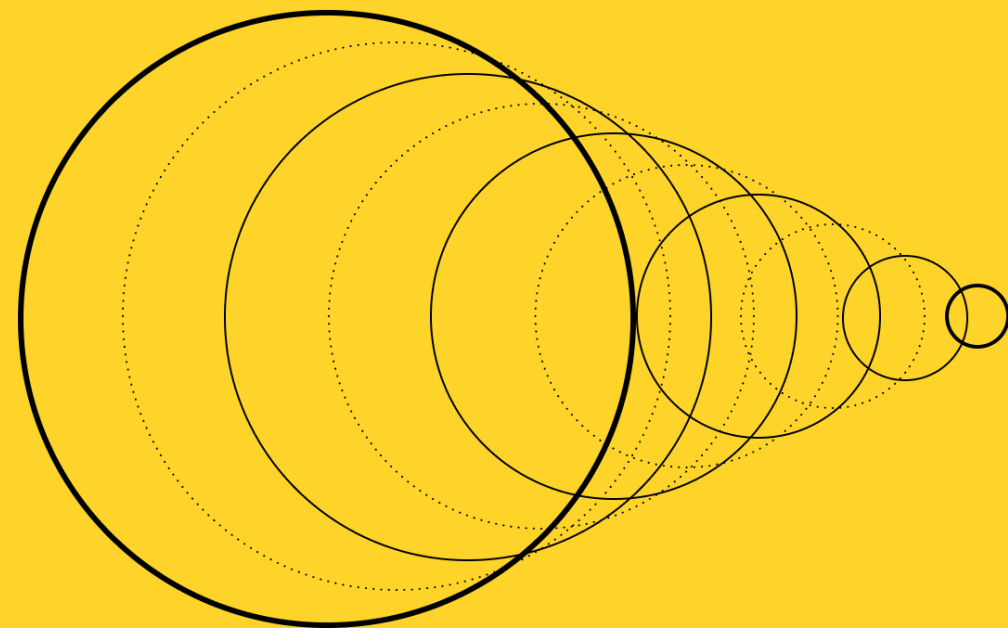* all conclusions are from our internal experiments

# DATA

# How to Distill

- Copy-paste from T2I does not (optimally) work

- Fun cross-task interactions:

    - LCM on editing > LCM on T2I for T2I quality 🥴

46

* all conclusions are from our internal experiments

# How to Distill

- For unmerged models, we can combine acceleration techniques
  - LoRA distill on the image part
  - Spec dec, KV-caching etc on the text part

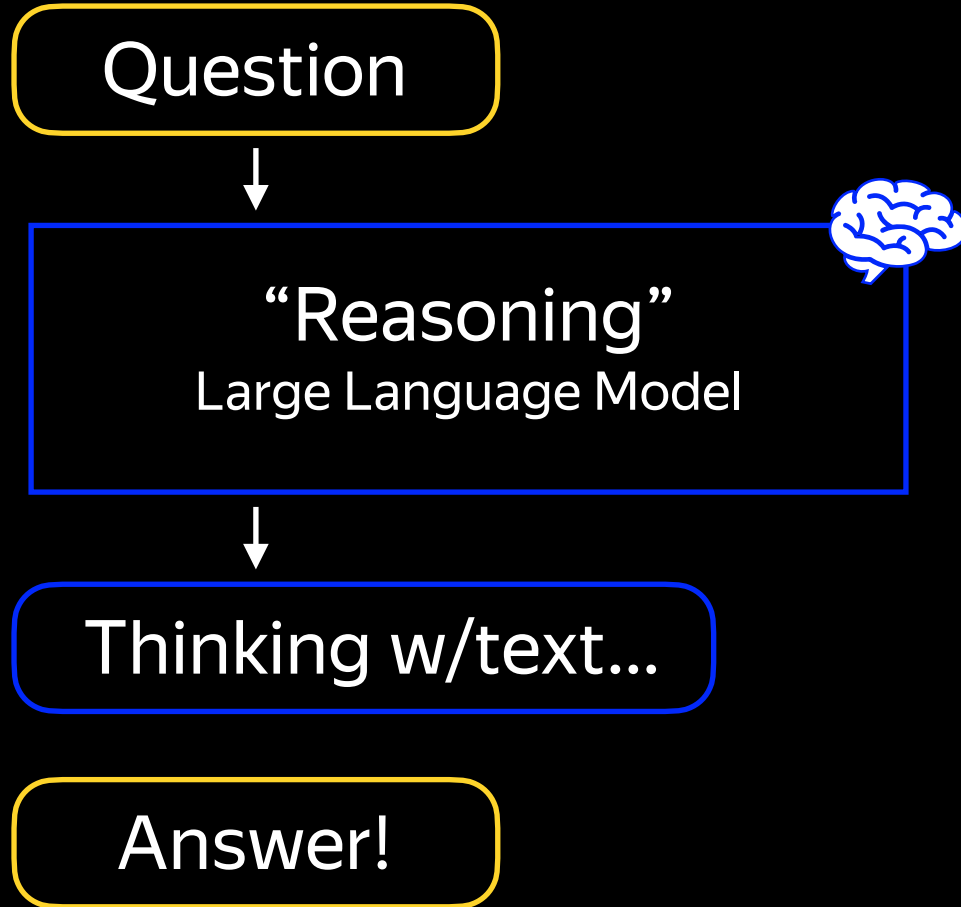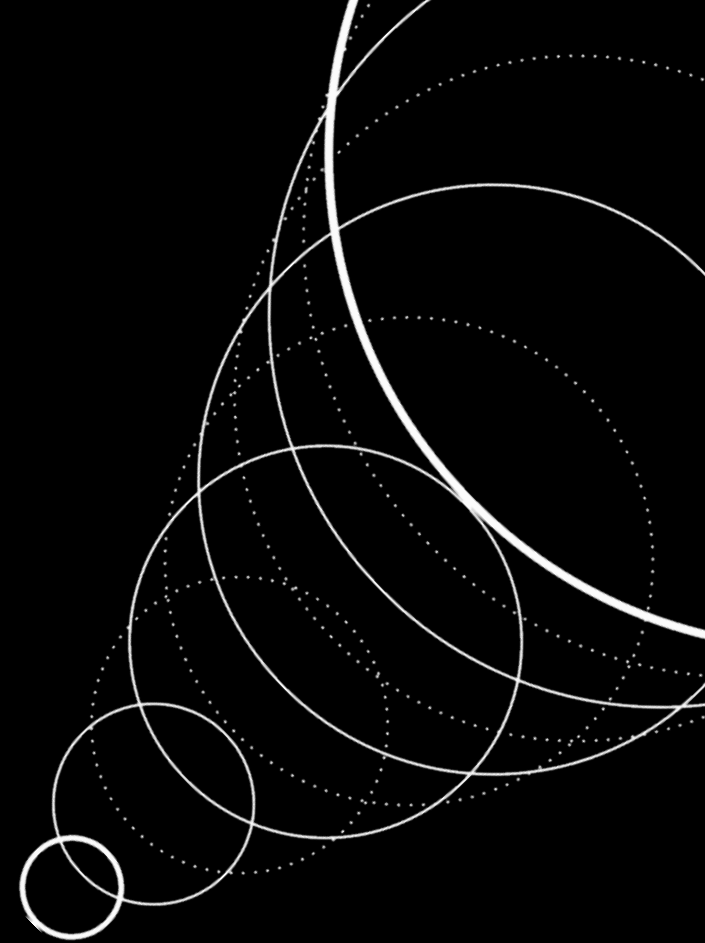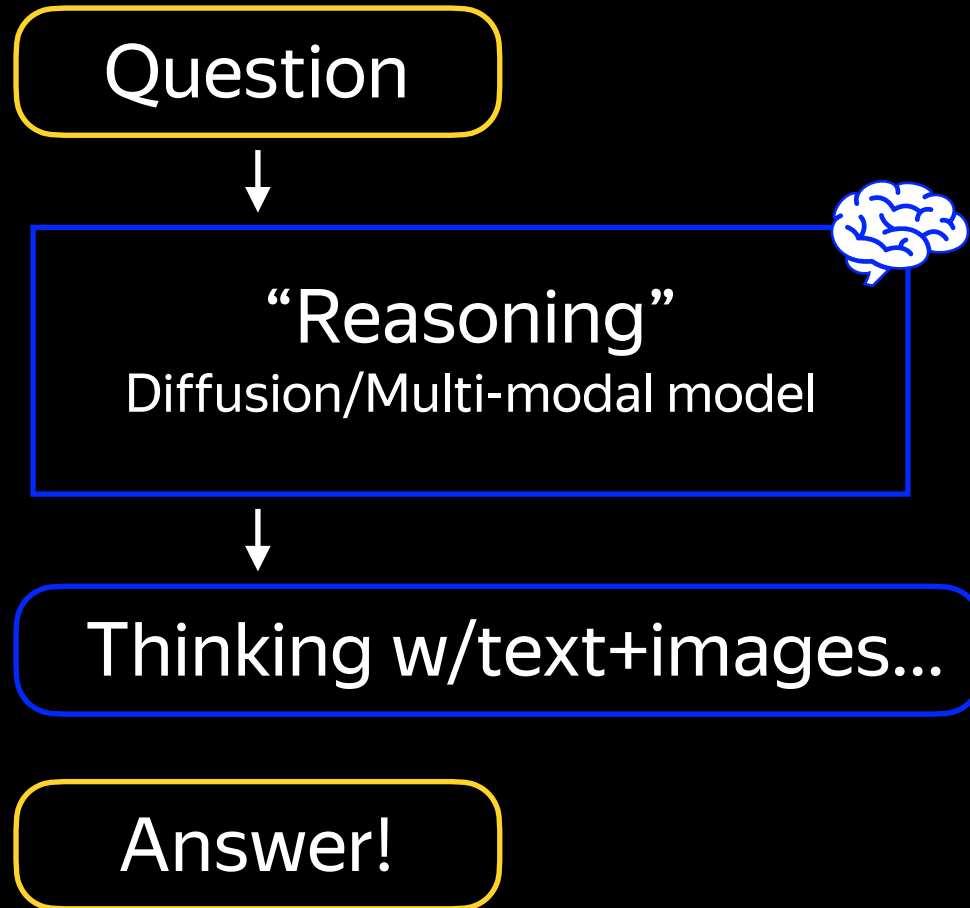- Can we use text-gen acceleration techniques on image-gen? 🤔
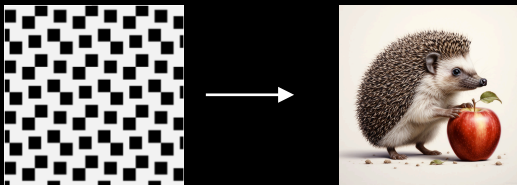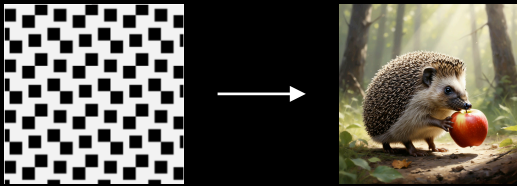
47

Yandex Research

# Inference-time Compute Scaling

# Very much known for LLMs



49

# Very much known for LLMs

Question

↓

"Reasoning"
Diffusion/Multi-modal model

↓

Thinking w/text+images...

Answer!

# Best-on-N baseline



"*hedgehog with an apple*"

E. Xie et. al., "SANA 1.5: Efficient Scaling of Training-Time and Inference-Time Compute in Linear Diffusion Transformer", 2025

# Best-on-N baseline



"*hedgehog with an apple*"

E. Xie et. al., "SANA 1.5: Efficient Scaling of Training-Time and Inference-Time Compute in Linear Diffusion Transformer", 2025

# Best-on-N baseline



"*hedgehog with an apple*"

E. Xie et. al., "SANA 1.5: Efficient Scaling of Training-Time and Inference-Time Compute in Linear Diffusion Transformer", 2025

# External VLM-as-judge

*Prompt* → T2I → *Image* → VLM → *Feedback* LLM → *New prompt*

M. Abdul et. al., "Test-time Prompt Refinement for Text-to-Image Models", 2025

# Can multi-modal model do all that itself?



*Prompt* → MLLM T2I → *Image* → MLLM VLM → *Feedback* → MLLM LLM → *New prompt*

# Boosts text gen and `many objects`

No TTS                    TTS



*"do an inscription `rash 435`"*

*"draw 5 dogs of different colors"*

* preliminary results

# Can we do the same for editing?



*Prompt + image* → **MLLM edit** → *Image* → **MLLM VLM** → *Feedback* → **MLLM LLM** → *New prompt*

# Also works pretty well

No TTS        TTS

*"change the action of the horses to galloping"*

*"Extract the person from the photo and dress them in a police uniform"*

* preliminary results

# Use all four tasks?

# Distill again?



**Prompt** → **MLLM T2I** → **Image** → **MLLM VLM** → **Feedback** → **MLLM LLM** → **New prompt**

**Prompt + image** → **MLLM edit** → **Image** → **MLLM VLM** → **Feedback** → **MLLM LLM** → **New prompt**
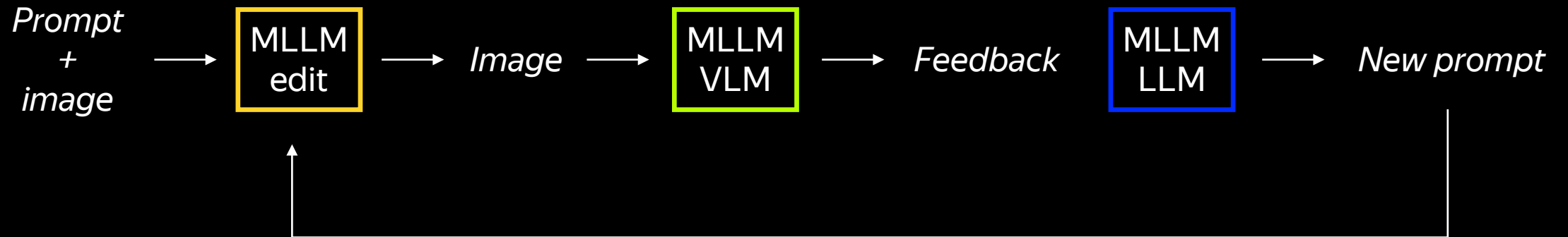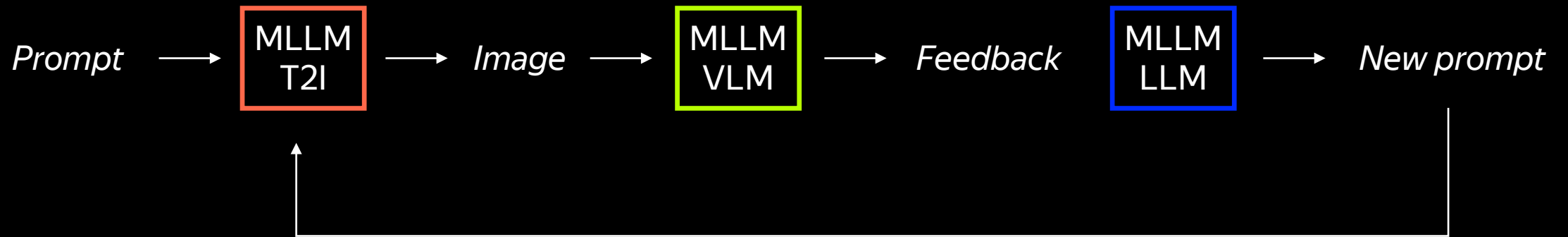
**What stops us from training a reasoner on this synthetic chains?**

# Instead of Conclusion

## Yandex Research also works on and presents:

**Poster 39**
Alchemist: Turning Public Text-to-Image Data into Generative Gold

**Poster 38**
Inverse Bridge Matching Distillation

**Poster 25**
TabM: Advancing tabular deep learning with parameter-efficient ensembling

**Poster 62**
Leveraging Coordinate Momentum in SignSGD and Muon: Memory-Optimized Zero-

**Poster 42**
AutoJudge: Judge Decoding Without Manual Annotation

**Poster 48**
Hogwild! Inference: Parallel LLM Generation via Concurrent Attention

**Poster 27**
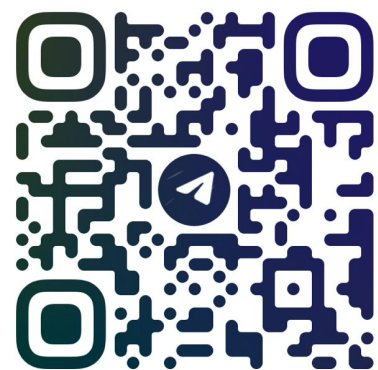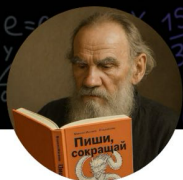GraphLand: Evaluating Graph Machine Learning Models on Diverse Industrial Data

**Poster 83**
On Linear Convergence in Smooth Convex-Concave Bilinearly-Coupled Saddle-Point Optimization: Lower Bounds and Optimal Algorithms

**Gen Models**         **Tabular DL**         **Effective Inference**         **Graph ML**         **Optimization**

# Yandex Research



Read research papers



Research with us



Develop Tech with us