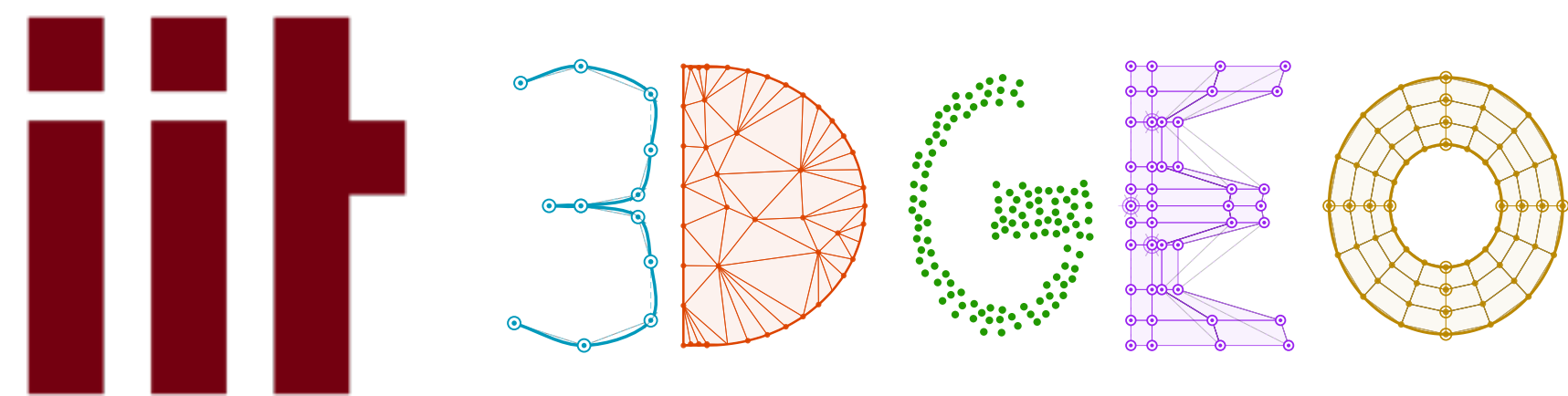
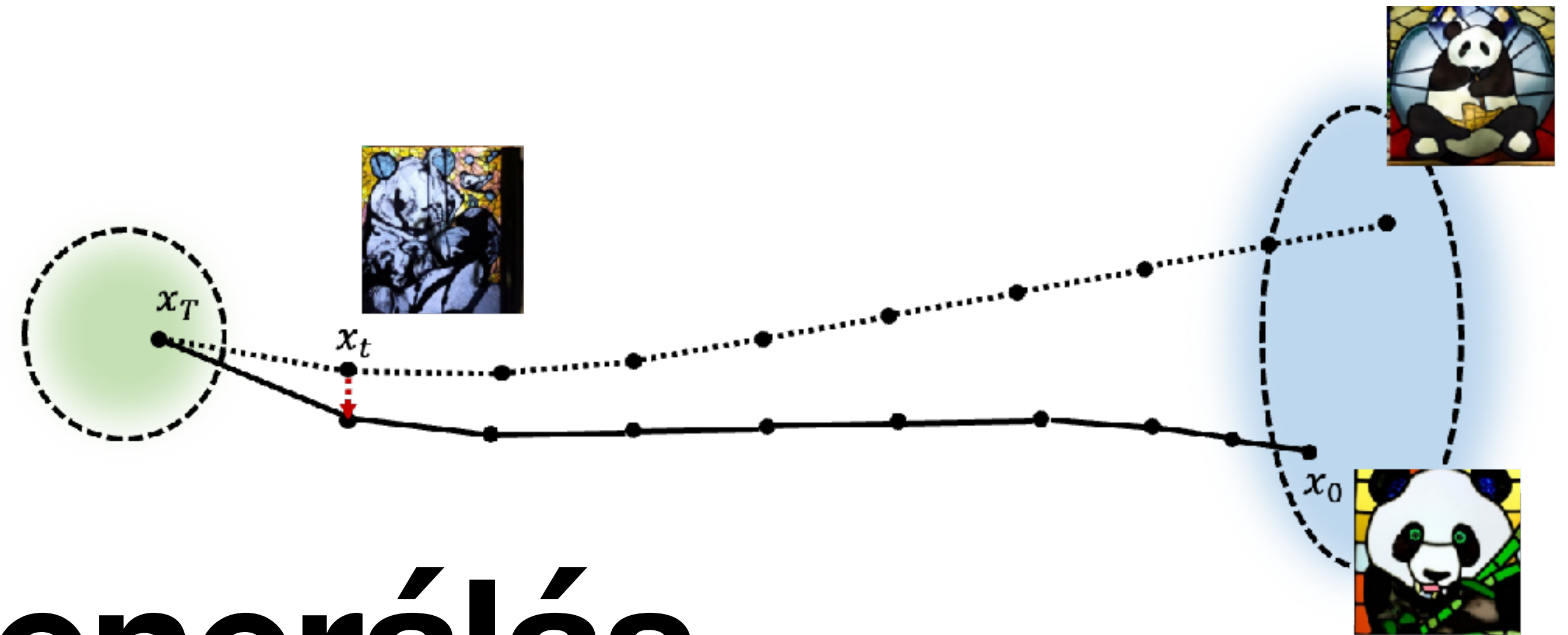


9. Előadás: Vezérelt Képgenerálás

Generatív AI és Inverz Módszerek a Képszintézisben
BME-VIK IIT, 2026

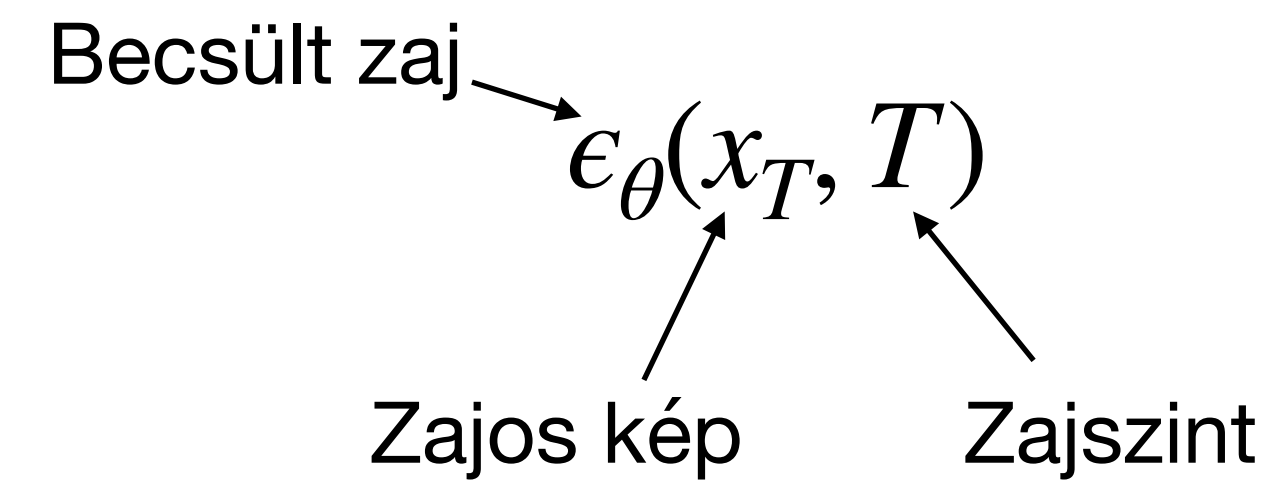


Dr. Vaitkus Márton

Vezérelt Generálás



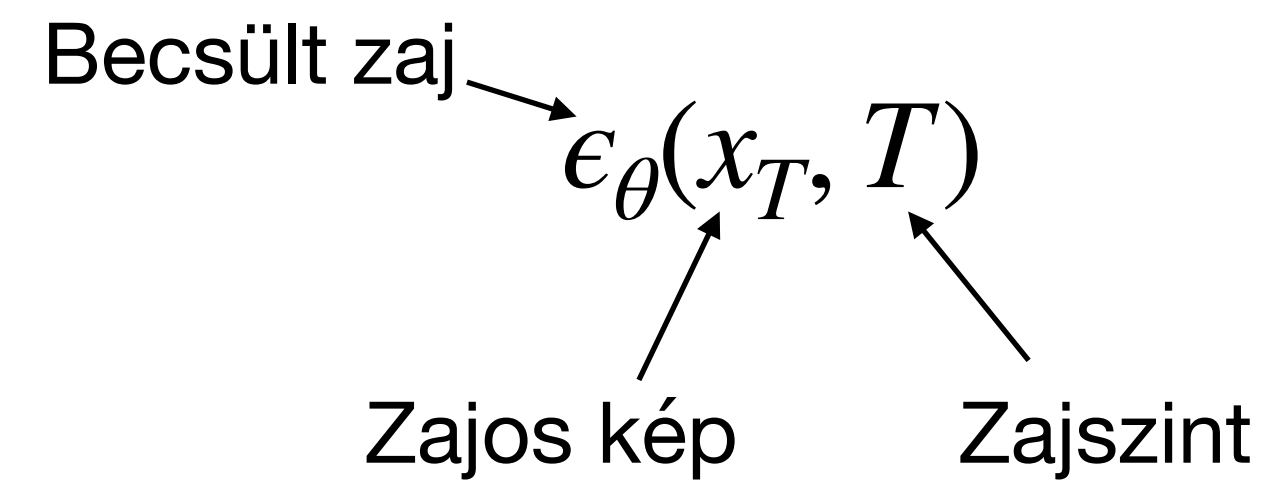
Diffúziós modell (*feltétel nélkül*):



Vezérelt Generálás



Diffúziós modell (*feltétel nélkül*):



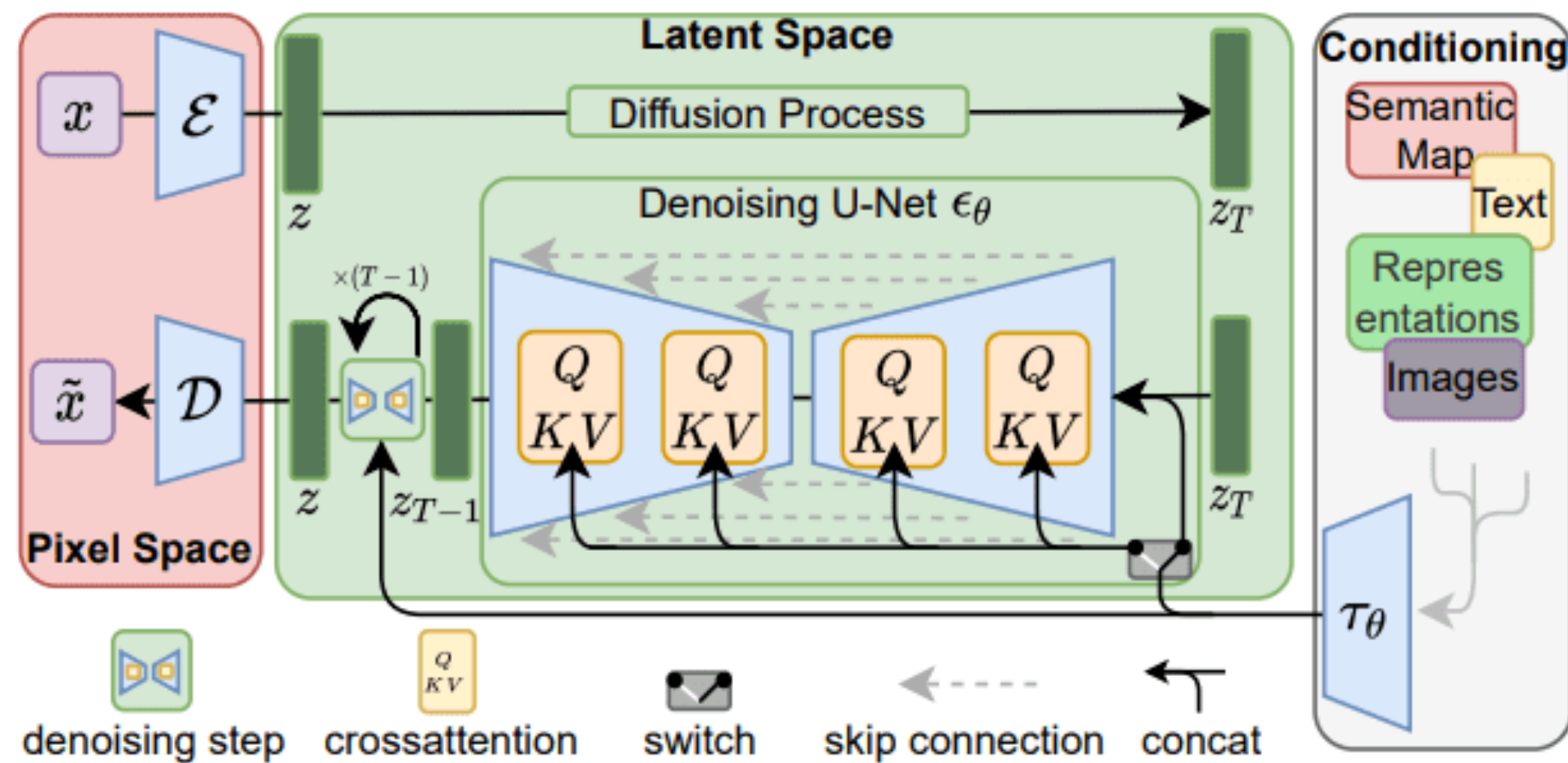
Diffúziós modell (*szöveges feltétellel*):

$$\epsilon_{\theta}(x_T, T, c)$$

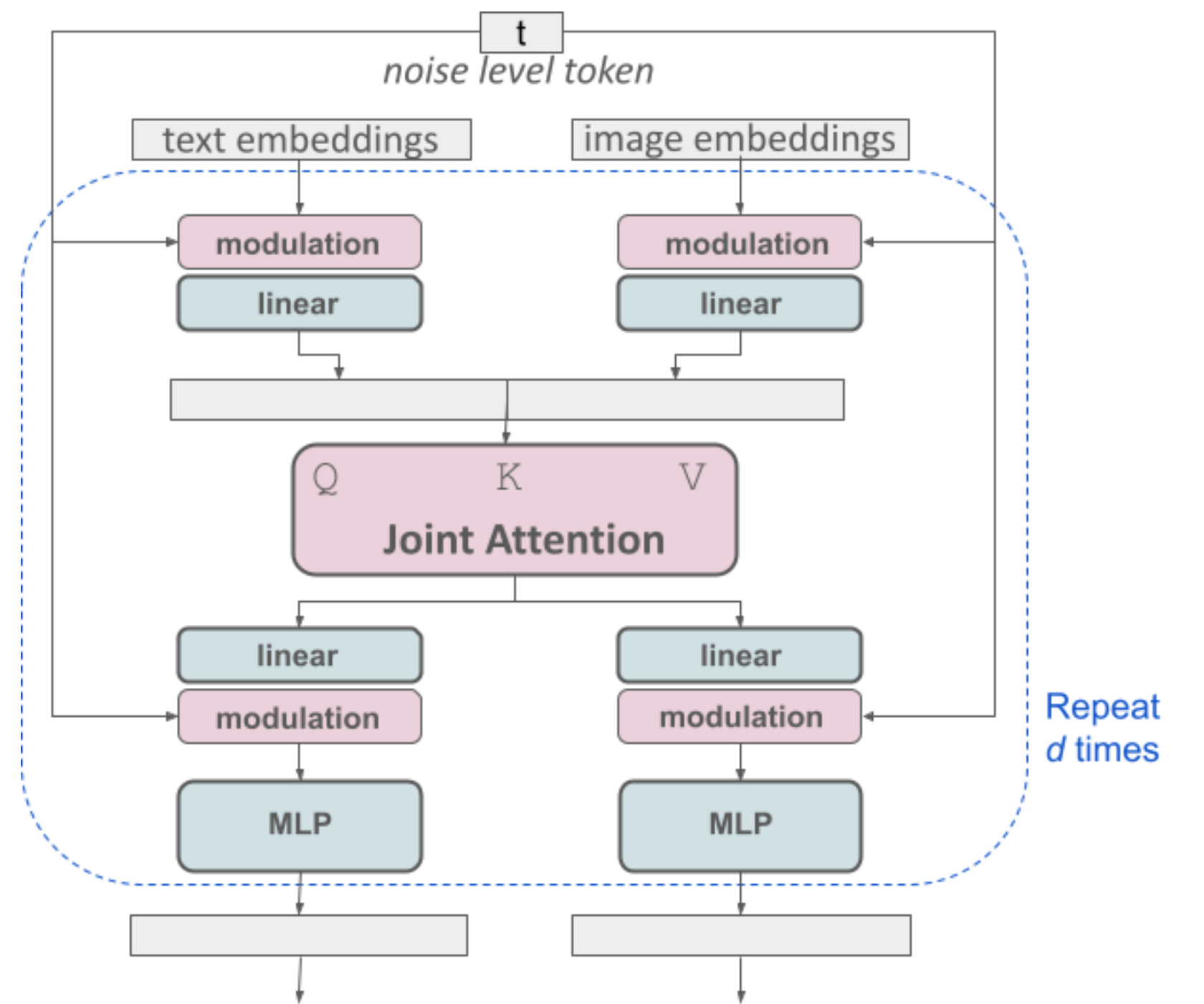
$c = \text{"A stained glass window of a panda eating bamboo."}$

Vezérelt Generálás

Neurális implementáció



Stable Diffusion 1.0
(U-Net + cross-attention)



Stable Diffusion 3.5
(Multimodális Transformer)

Vezérelt Generálás

Instrukciók követése

“A stained glass window of a panda eating bamboo.”

Vezérelt Generálás

Instrukciók követése

“A stained glass window of a panda eating bamboo.”



Vezérelt Generálás

Instrukciók követése

“A stained glass window of a panda eating bamboo.”



Tudunk jobbat?

Vezérelt Generálás

Instrukciók követése

“A stained glass window of a panda eating bamboo.”



Tudunk jobbat?



Vezérelt Generálás

Instrukciók követése

“A stained glass window of a panda eating bamboo.”



Tudunk jobbat?

Igen!

“Guidance”



Vezérelt Generálás

Instrukciók követése

“A stained glass window of a panda eating bamboo.”



Tudunk jobbat?

Igen!

“Guidance”

*Kritikus eleme
minden
képgenerátornak!!!*



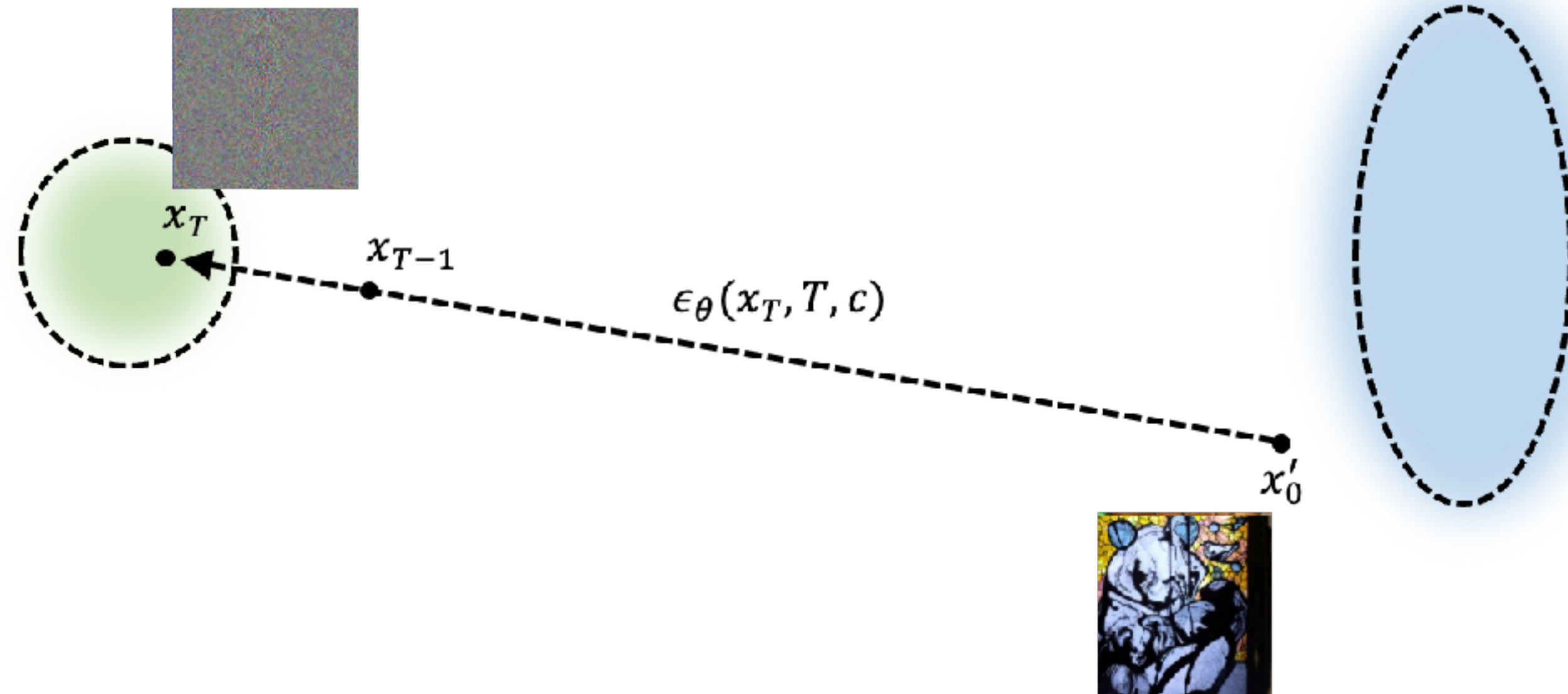
Vezérelt Generálás

Classifier Guidance

Vezérelt Generálás

Classifier Guidance

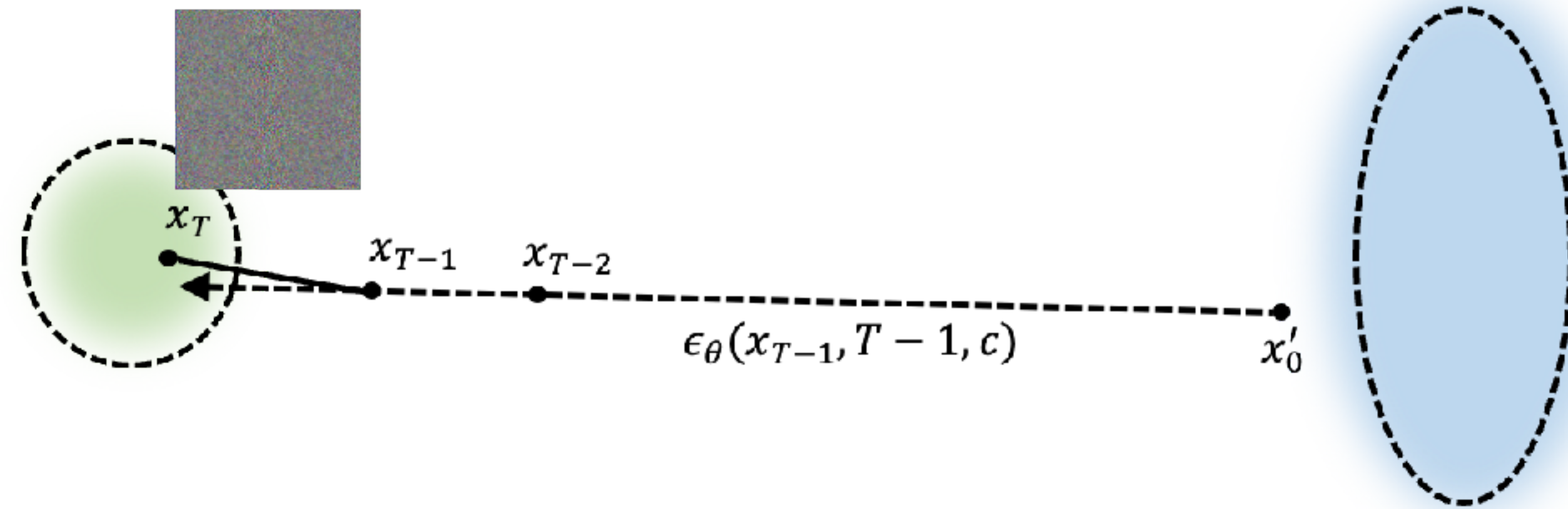
$c = \text{"A stained glass window of a panda eating bamboo."}$



Vezérelt Generálás

Classifier Guidance

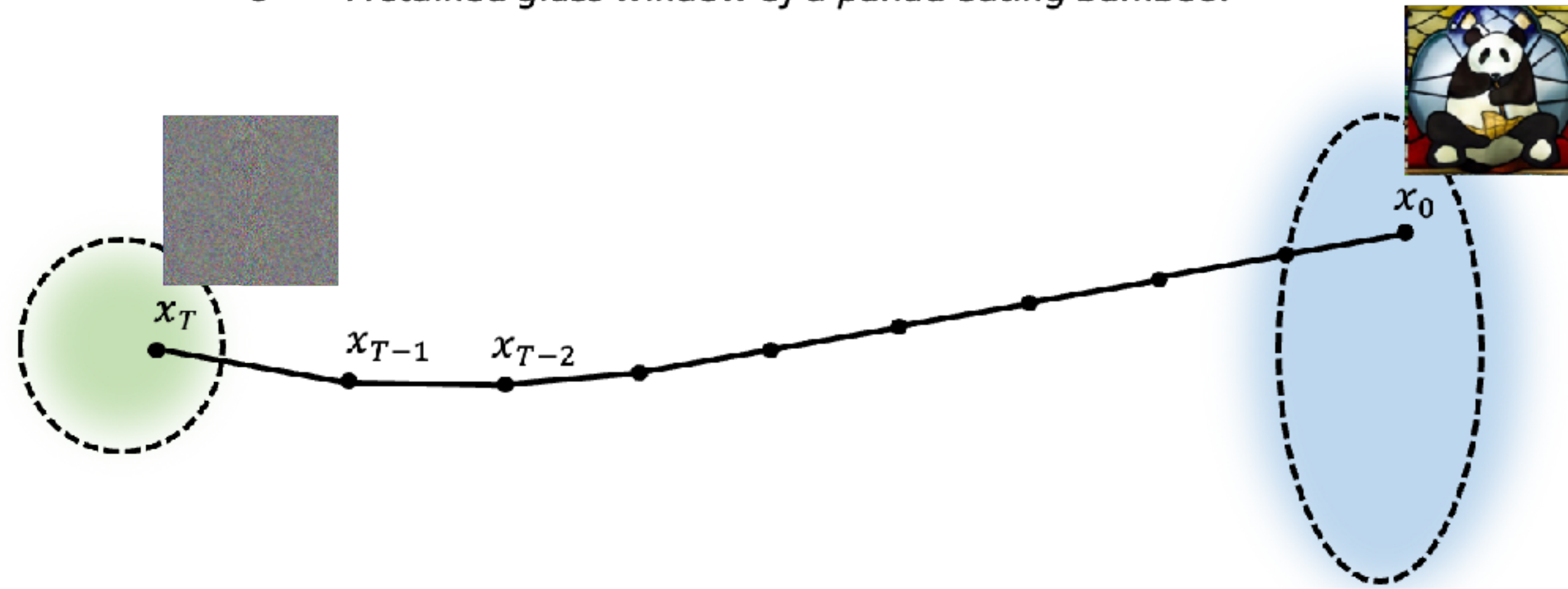
$c = \text{"A stained glass window of a panda eating bamboo."}$



Vezérelt Generálás

Classifier Guidance

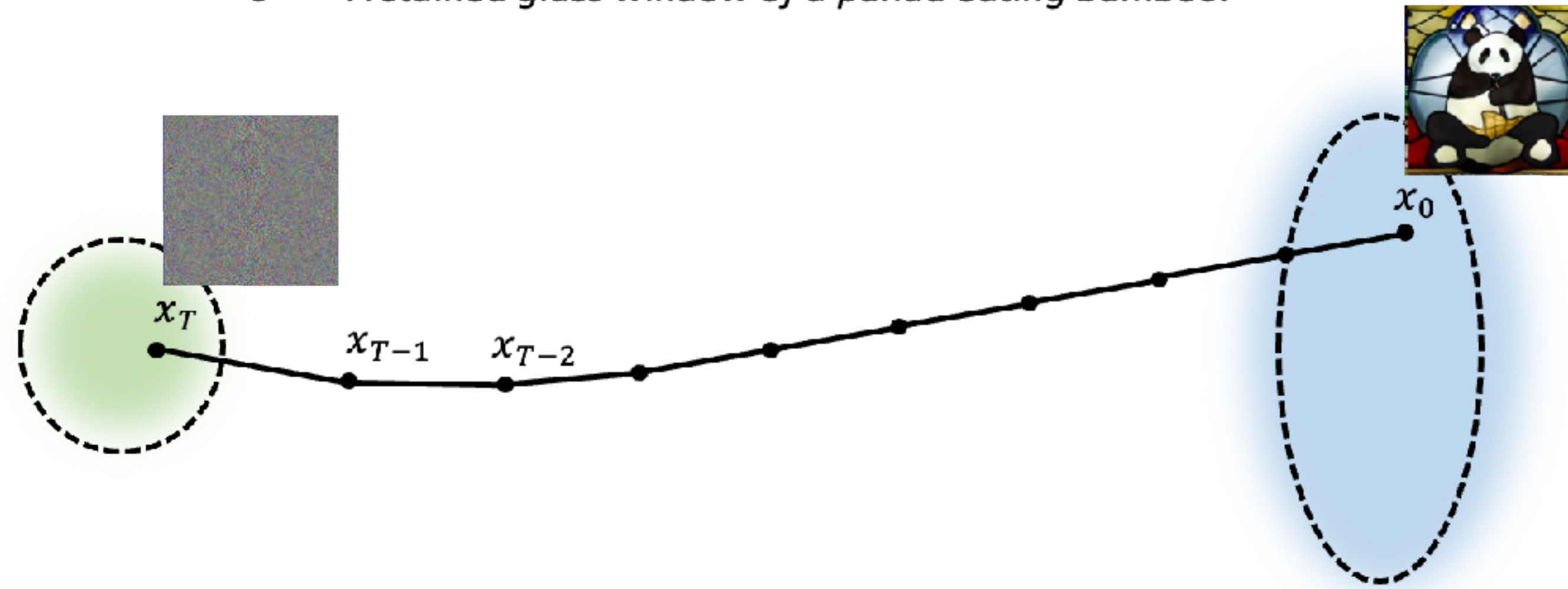
$c = \text{"A stained glass window of a panda eating bamboo."}$



Vezérelt Generálás

Classifier Guidance

$c = \text{"A stained glass window of a panda eating bamboo."}$

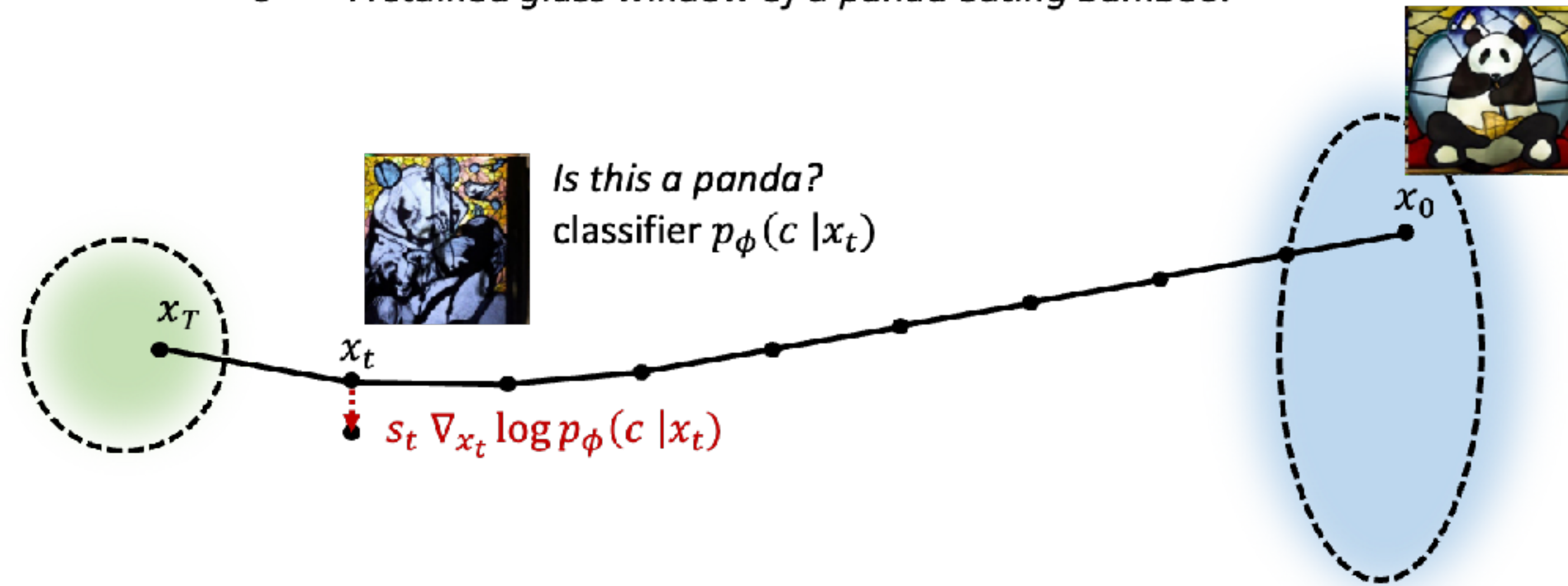


Classifier Guidance — minősítjük a zajos kép prompthoz való igazodását — $p_\phi(c | x_t)$

Vezérelt Generálás

Classifier Guidance

$c = \text{"A stained glass window of a panda eating bamboo."}$



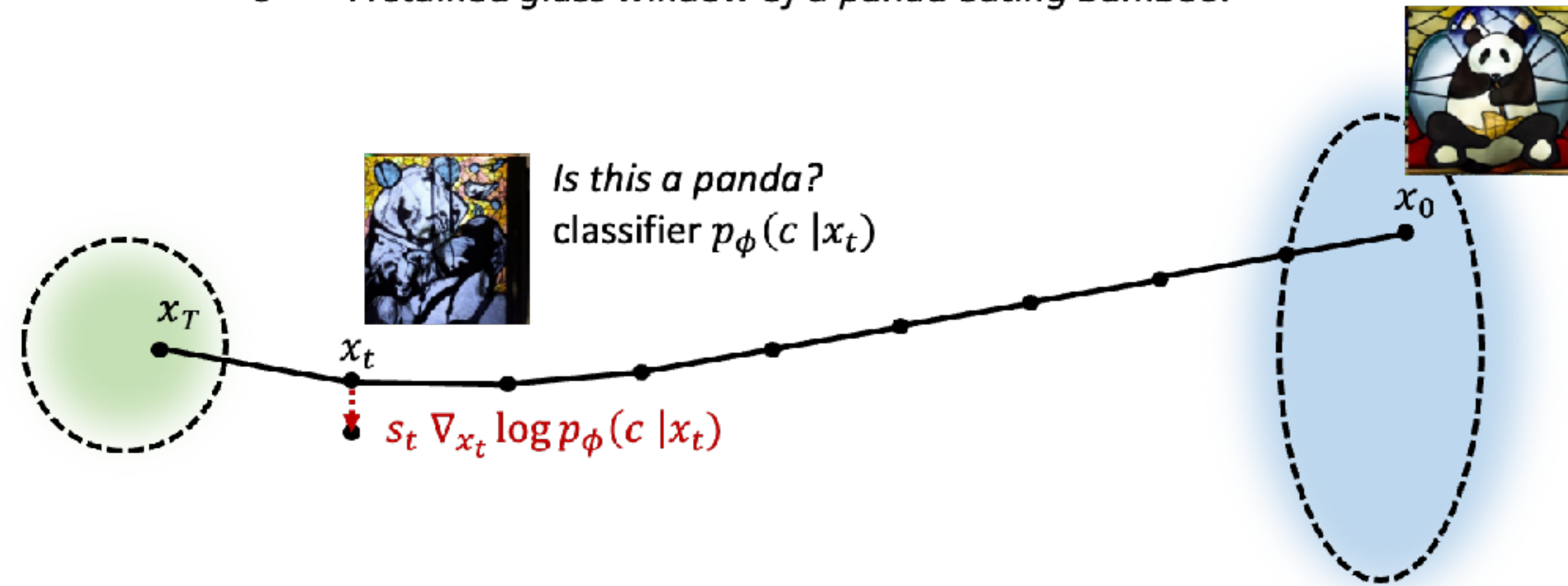
Classifier Guidance — minősítsük a zajos kép prompthoz való igazodását — $p_\phi(c | x_t)$

Bayes tételből: $\nabla_x \log p(x_t | c) = \nabla_x \log p(c | x_t) + \nabla_x \log p(x_t) - \nabla_x \log p(c)$

Vezérelt Generálás

Classifier Guidance

$c = \text{"A stained glass window of a panda eating bamboo."}$



Classifier Guidance — minősítsük a zajos kép prompthoz való igazodását — $p_\phi(c | x_t)$

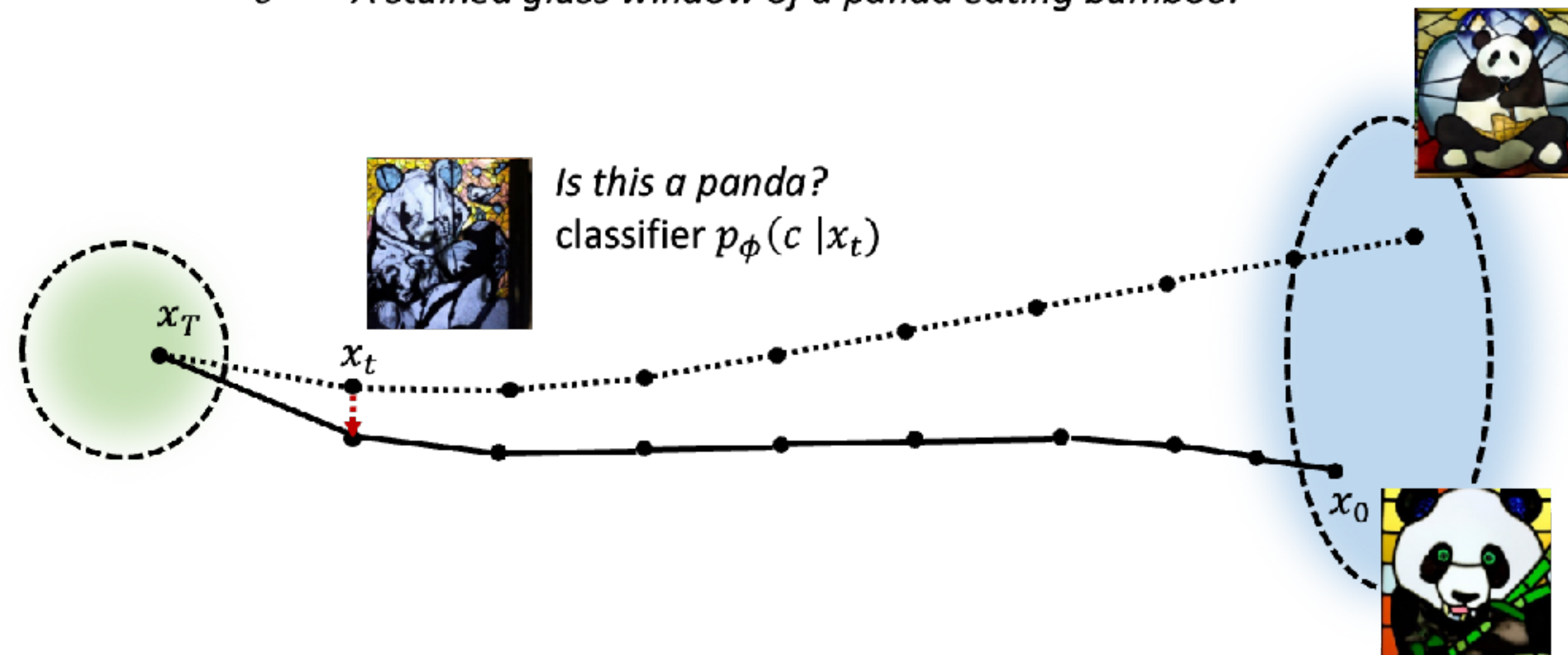
Bayes tételből: $\nabla_x \log p(x_t | c) = \nabla_x \log p(c | x_t) + \nabla_x \log p(x_t) - \nabla_x \log p(c)$

Skálázzuk fel a klasszifikáció erejét: $\nabla_x \log p(x_t) + s_t \cdot \nabla_x \log p(c | x_t)$

Vezérelt Generálás

Classifier Guidance

$c = \text{"A stained glass window of a panda eating bamboo."}$



Classifier Guidance — minősítsük a zajos kép prompthoz való igazodását — $p_\phi(c | x_t)$

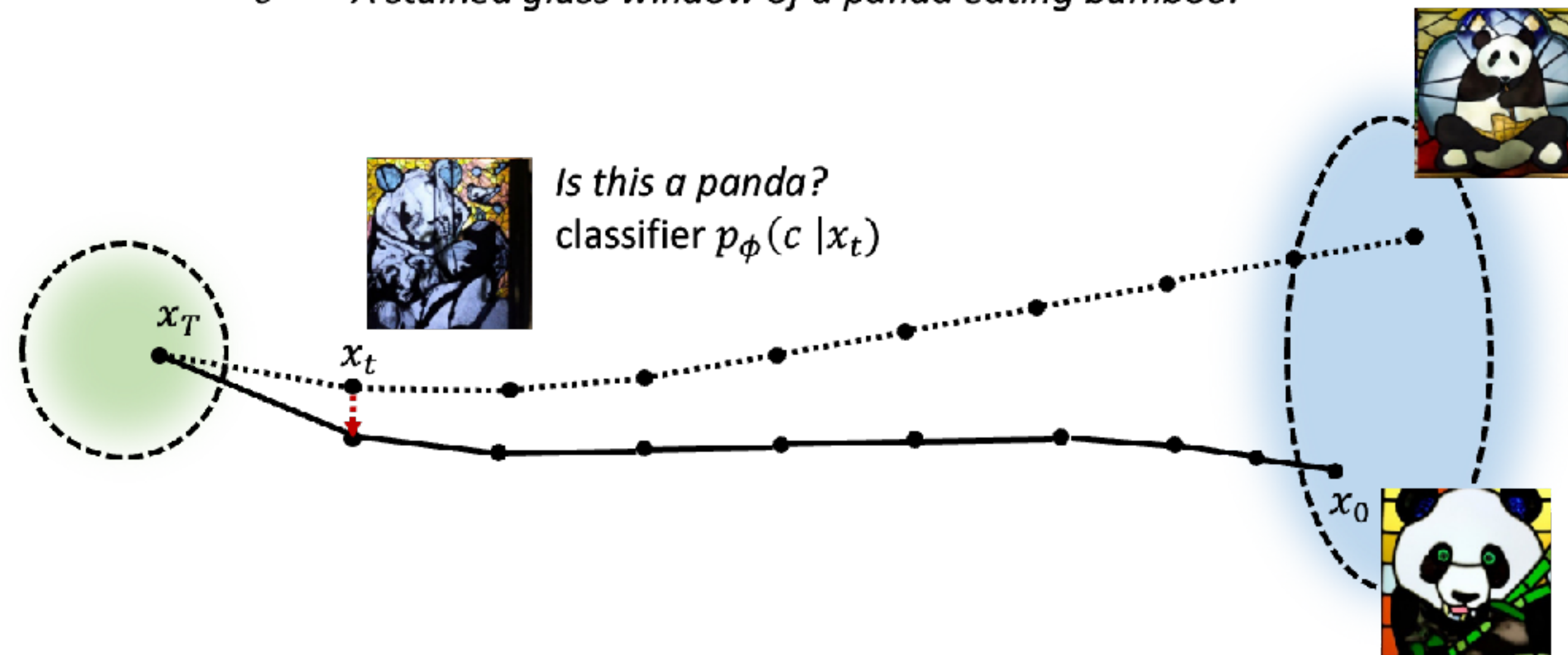
Bayes tételből: $\nabla_x \log p(x_t | c) = \nabla_x \log p(c | x_t) + \nabla_x \log p(x_t) - \nabla_x \log p(c)$

Skálázzuk fel a klasszifikáció erejét: $\nabla_x \log p(x_t) + s_t \cdot \nabla_x \log p(c | x_t)$

Vezérelt Generálás

Classifier Guidance

$c = \text{"A stained glass window of a panda eating bamboo."}$



Classifier Guidance — minősítsük a zajos kép prompthoz való igazodását — $p_\phi(c | x_t)$

Bayes tételből: $\nabla_x \log p(x_t | c) = \nabla_x \log p(c | x_t) + \nabla_x \log p(x_t) - \nabla_x \log p(c)$

Skálázzuk fel a klasszifikáció erejét: $\nabla_x \log p(x_t) + s_t \cdot \nabla_x \log p(c | x_t)$

Ehhez külön klasszifikátort kell tanítani, még hozzá zajos képeken!

Vezérelt Generálás

Classifier-Free Guidance (CFG)

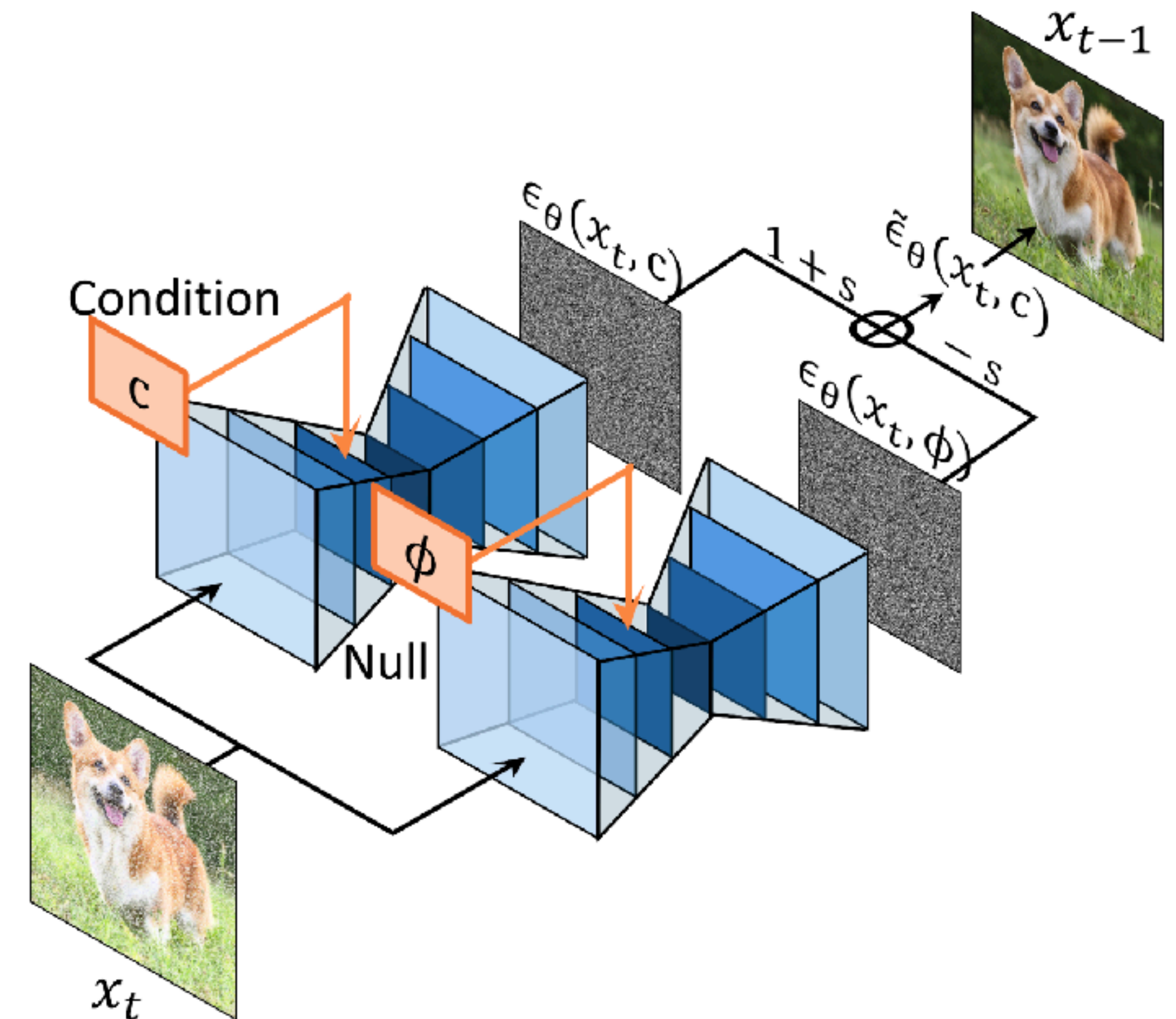
- Szabaduljunk meg a klasszifikációtól (Bayes újból):

$$\begin{aligned}
 & \nabla_x \log p(x_t) + s_t \cdot \nabla_x \log p(y | x_t) \\
 &= \nabla_x \log(p(x_t)) + s_t \cdot (\nabla_x \log p(x_t | y) + \cancel{\nabla_x \log p(y)} - \nabla_x \log p(x_t)) \\
 &= (1 - s_t) \underbrace{\nabla_x \log p(x_t)}_{\text{Denoizer feltétel nélkül}} + s_t \cdot \underbrace{\nabla_x \log p(x_t | y)}_{\text{Denoizer feltétellel}}
 \end{aligned}$$

- Classifier-Free Guidance (CFG)** — a denoizert 2x értékeljük ki: feltétellel és anélkül (“null” feltétel), majd kombináljuk!
- Címkézett képekkel való tanítás során valamilyen valószínűséggel elhagyjuk a feltételt (“null”-al helyettesítjük), hogy a háló feltétel nélkül is megtanuljon generálni!

CLASSIFIER-FREE DIFFUSION GUIDANCE

Jonathan Ho & Tim Salimans
Google Research, Brain team

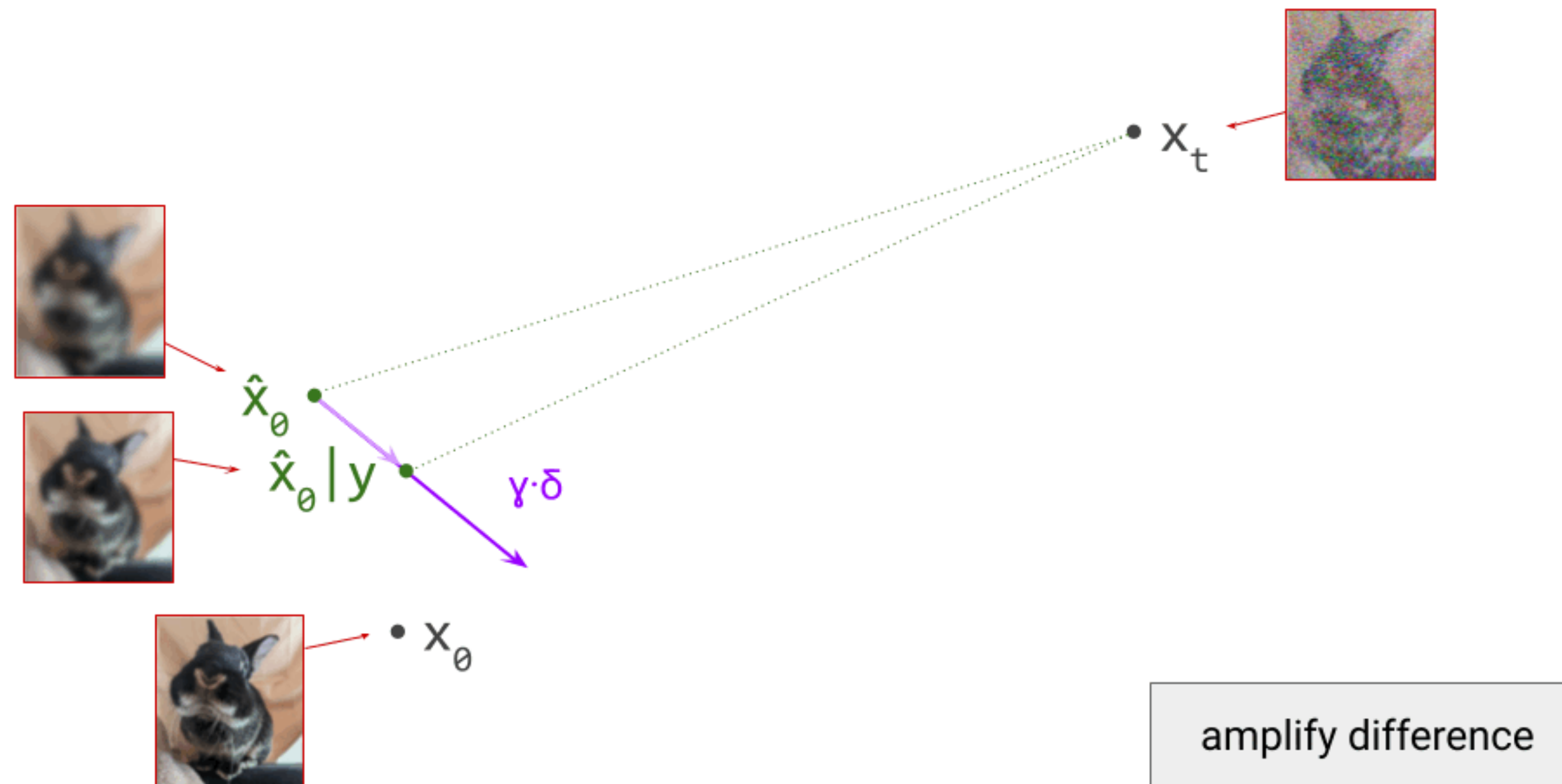


Vezérelt Generálás

Classifier-Free Guidance (CFG)

CLASSIFIER-FREE DIFFUSION GUIDANCE

Jonathan Ho & Tim Salimans
Google Research, Brain team



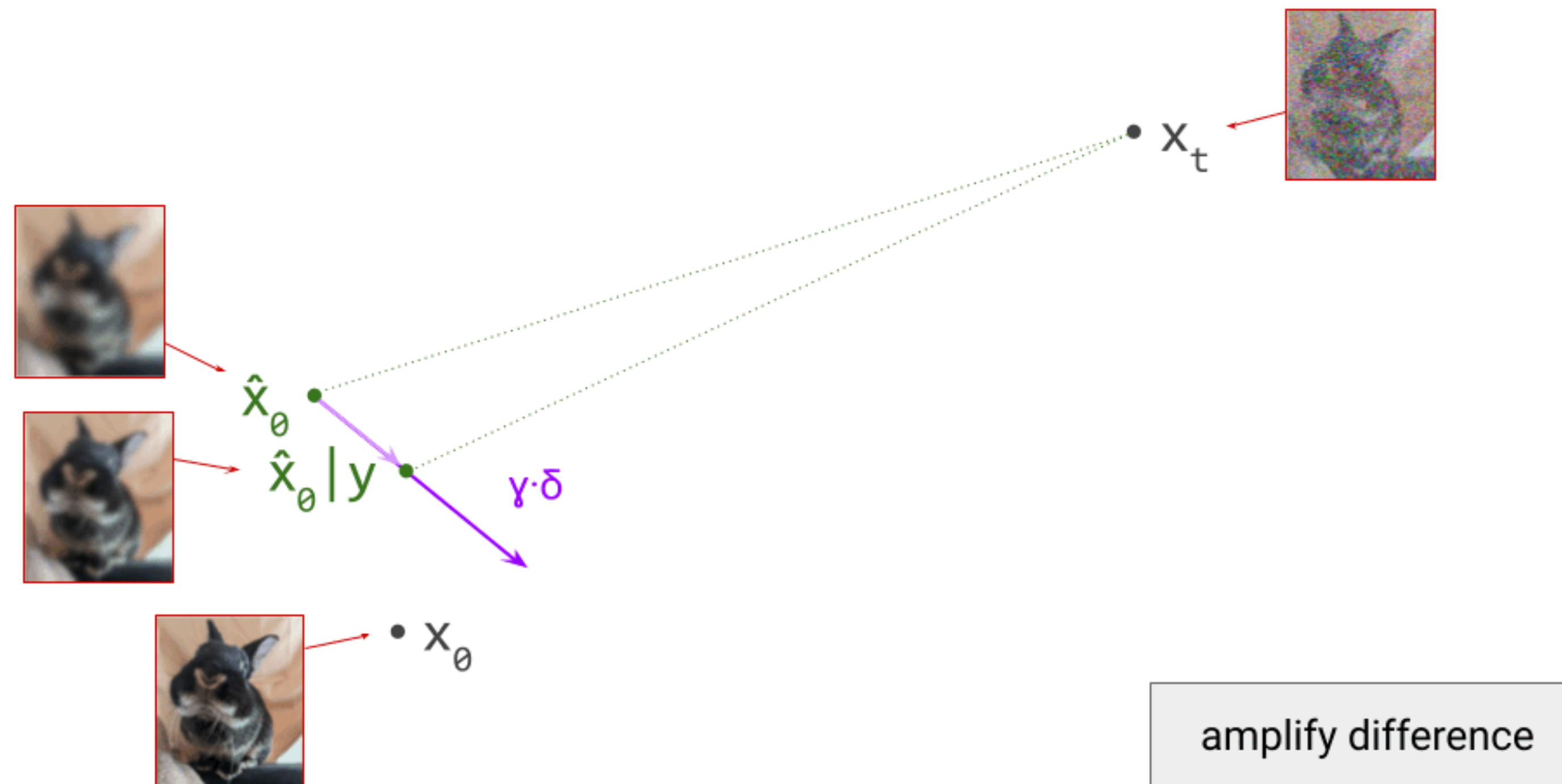
DDPM + CFG — Forrás: [LINK](#)

Vezérelt Generálás

Classifier-Free Guidance (CFG)

CLASSIFIER-FREE DIFFUSION GUIDANCE

Jonathan Ho & Tim Salimans
Google Research, Brain team



DDPM + CFG — Forrás: [LINK](#)

Vezérelt Generálás

Classifier-Free Guidance (CFG)

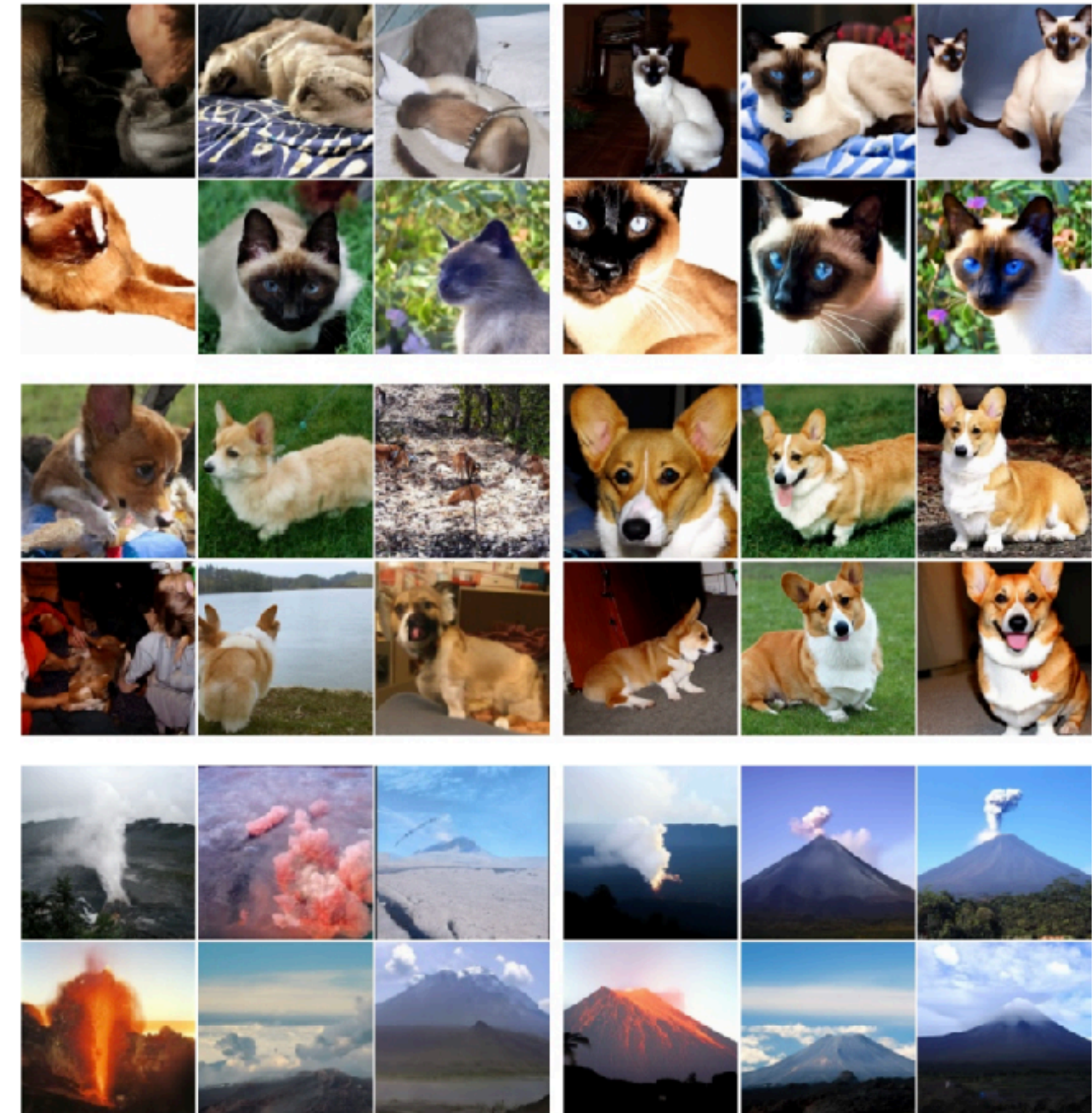
- Tapasztalati tény: erősebb CFG skálafaktorral ($s_t > 1$) jobb eredményeket kapunk!!!
- Gyakorlatban akár $s_t = 5 - 8$!!!
- Minden “rendes” képgenerátor erős CFG-t használ!

$s_t = 0$ 



CLASSIFIER-FREE DIFFUSION GUIDANCE

Jonathan Ho & Tim Salimans
Google Research, Brain team



$s_t = 0$

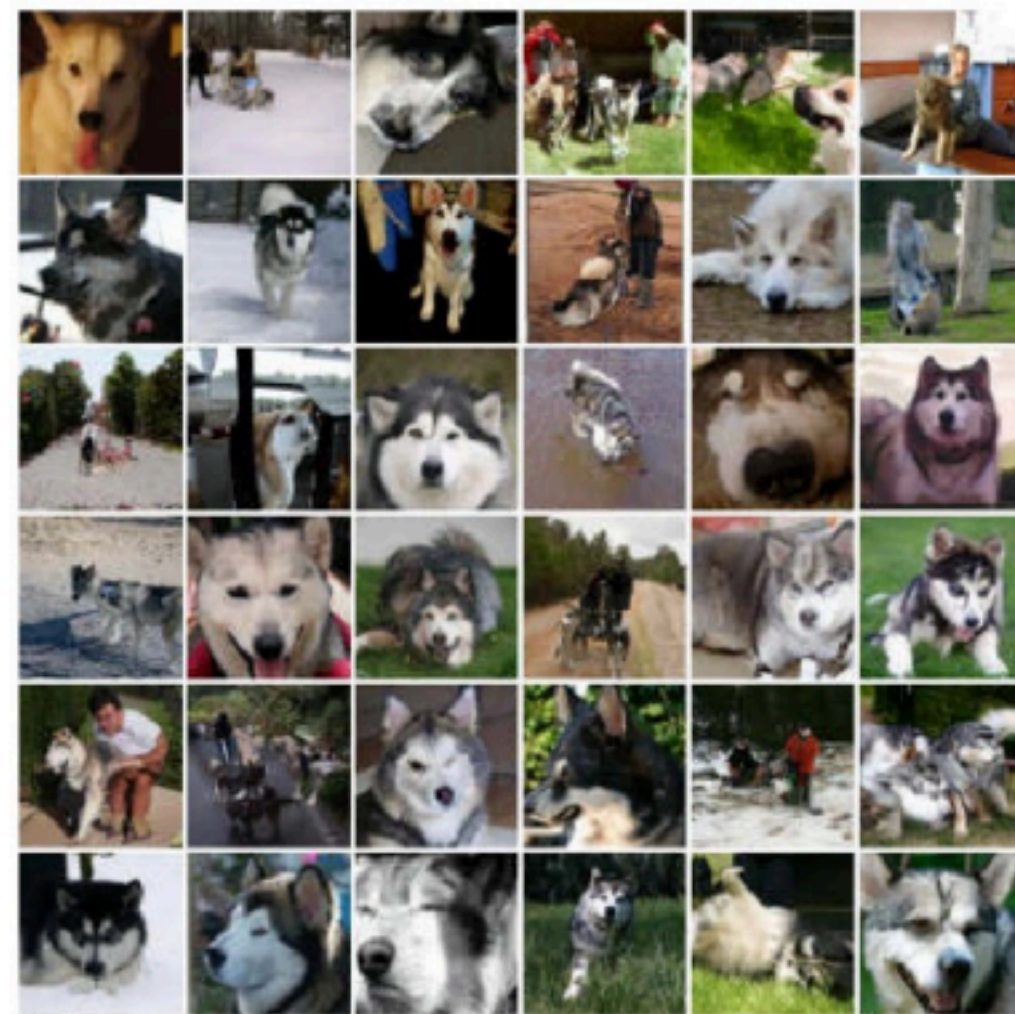
$s_t = 3$

Vezérelt Generálás

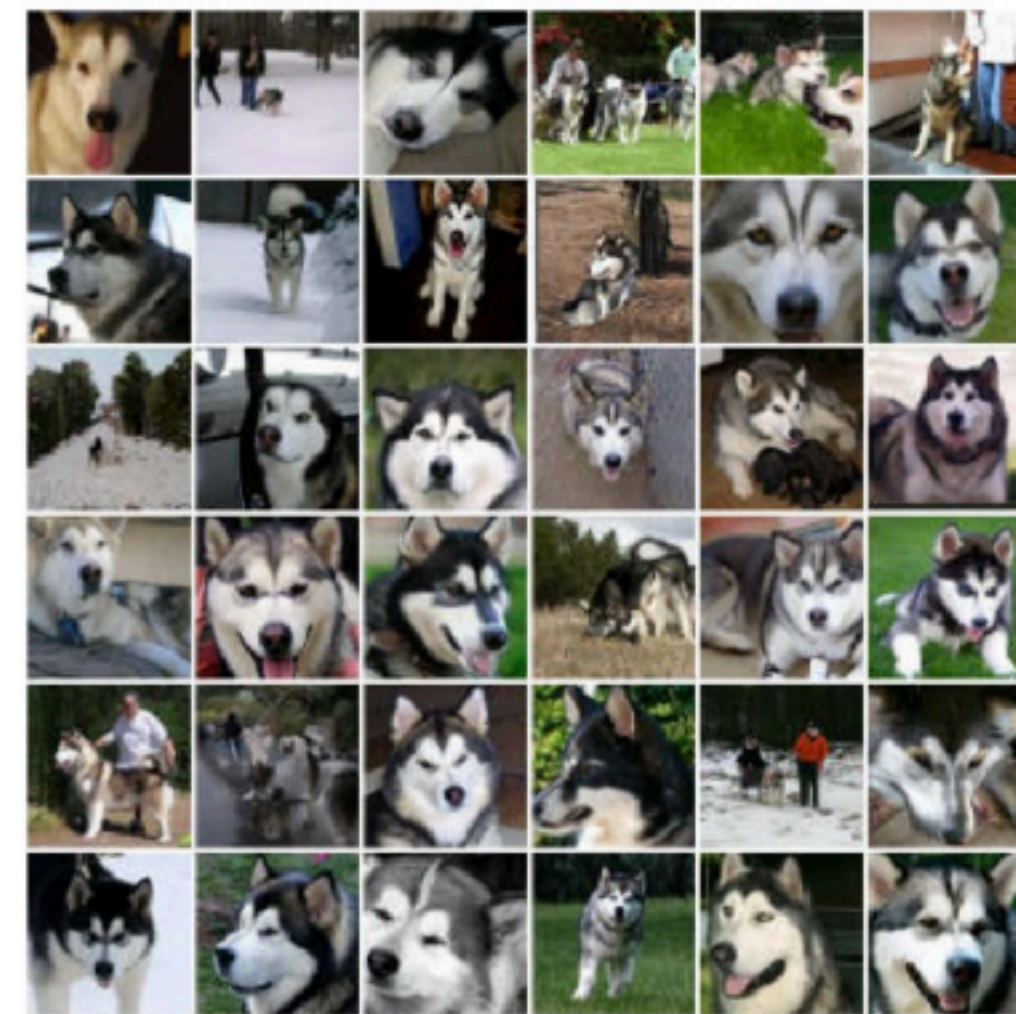
Classifier-Free Guidance (CFG)

CLASSIFIER-FREE DIFFUSION GUIDANCE

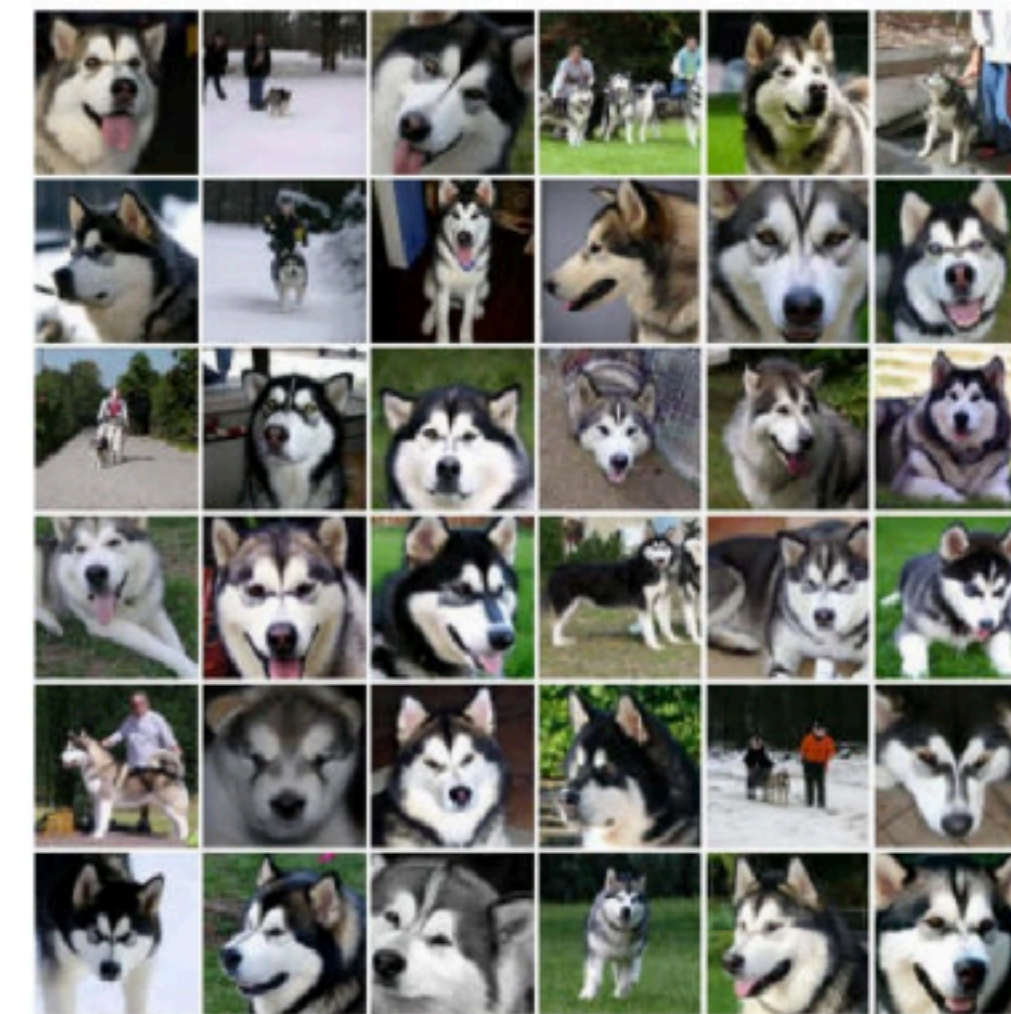
Jonathan Ho & Tim Salimans
Google Research, Brain team



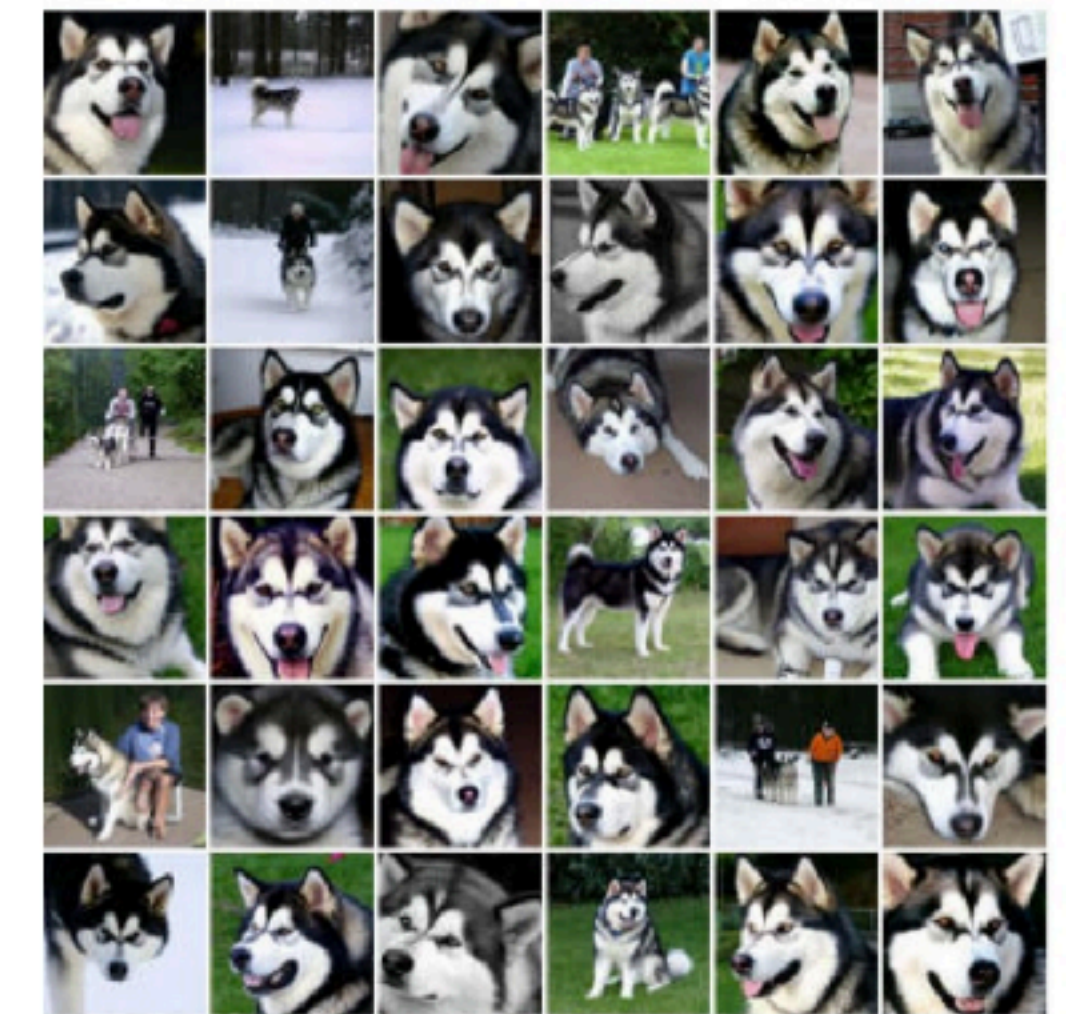
(a) $w = 0$



(b) $w = 1$



(c) $w = 2$



(d) $w = 4$

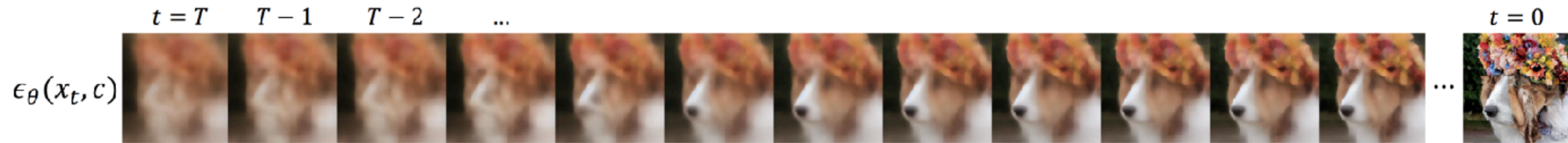
CFG skálázás hatása: jobb igazodás a prompthoz, de kisebb variancia és “szaturált” színek!

Vezérelt Generálás

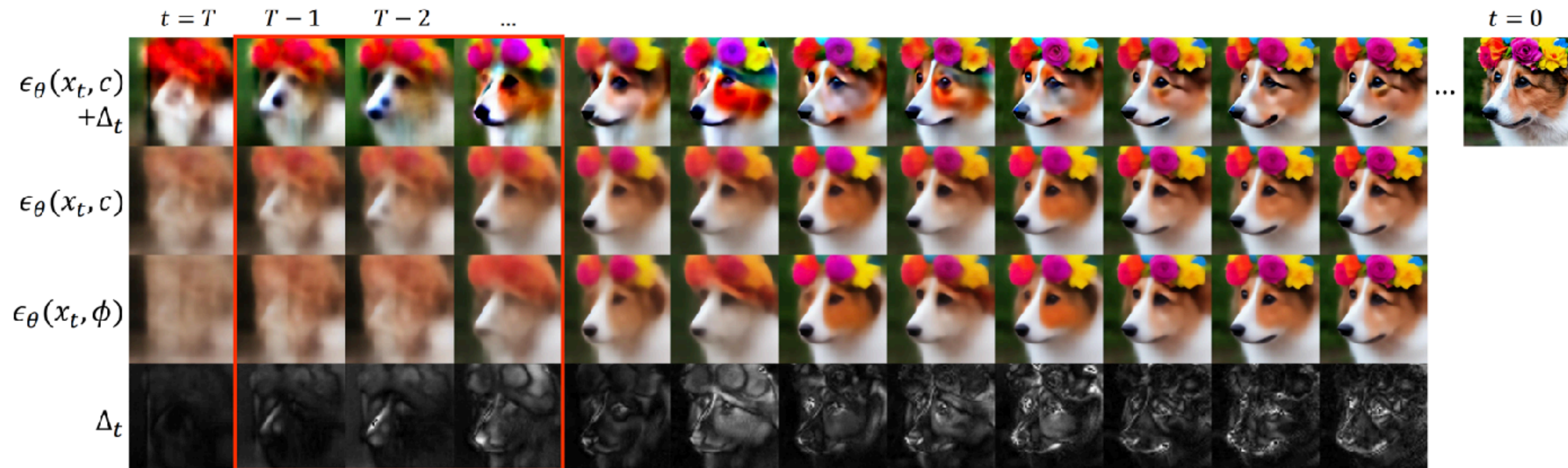
Classifier-Free Guidance (CFG)

CLASSIFIER-FREE DIFFUSION GUIDANCE

Jonathan Ho & Tim Salimans
Google Research, Brain team



(a) Diffusion sampling without CFG



(b) Diffusion sampling with CFG

A CFG hamarabb a promptnak megfelelő irányba “húzza” a generálást!

Vezérelt Generálás

Perturbed-Attention Guidance (PAG)

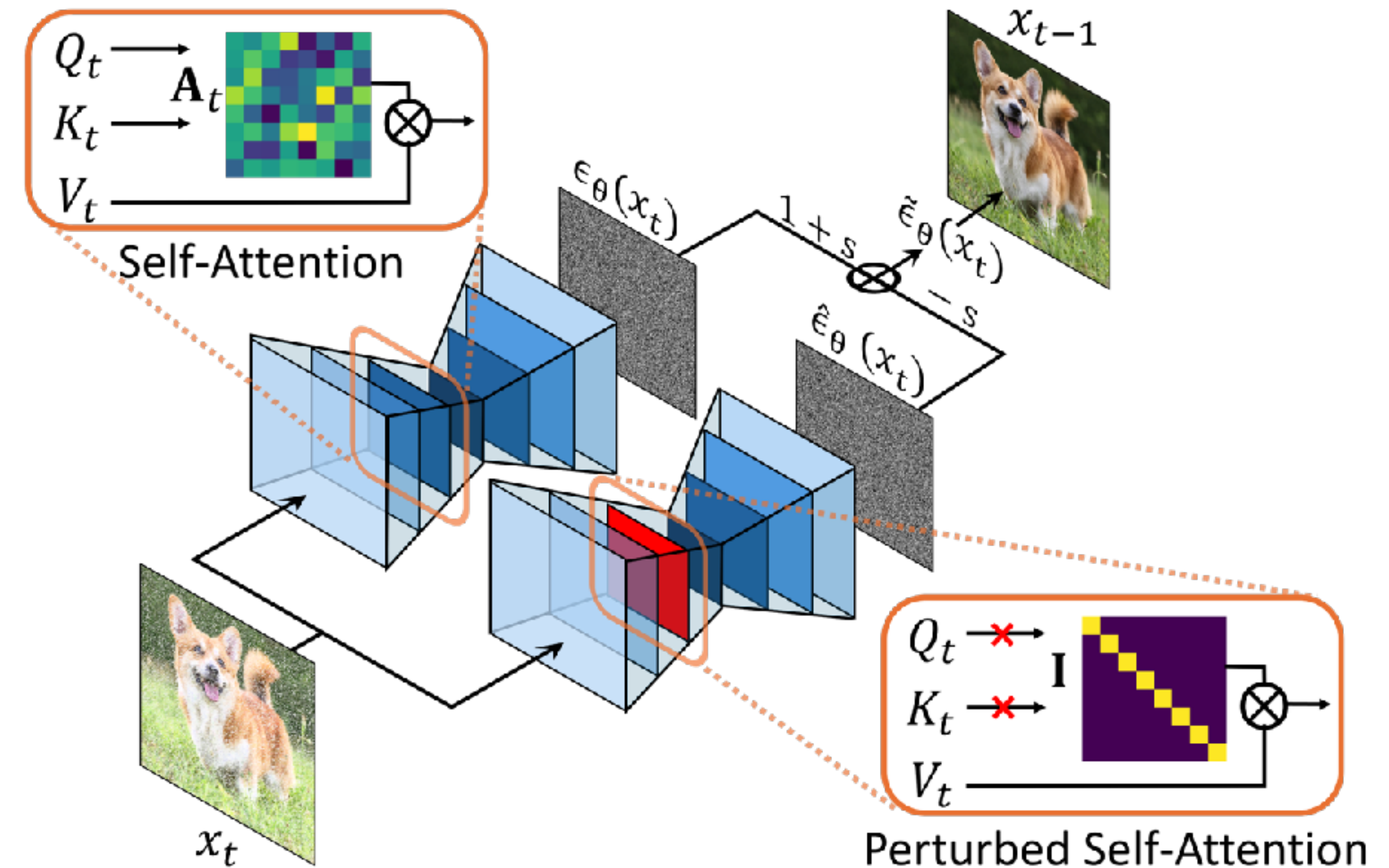
- Megfigyelés: a generálás korai (zajosabb) szakaszában néha “rossz” struktúrájú képek felé indulunk, amit nehéz korrigálni...
- A struktúrát a self-attention (query-key szorzat) mátrixok határozzák meg!
- **Perturbed-Attention Guidance (PAG):** CFG-hez hasonló, csak feltétel nélküli kiértékelés helyett a self-attention mátrixokat “perturbáljuk” (egységmátrixra cseréljük) és kombináljuk a normál attention eredményével!

Self-Rectifying Diffusion Sampling with Perturbed-Attention Guidance

Donghoon Ahn⁺¹, Hyoungwon Cho⁺¹, Jaewon Min¹, Wooseok Jang¹,
Jungwoo Kim¹, SeonHwa Kim¹, Hyun Hee Park²,
Kyong Hwan Jin¹, and Seungryong Kim¹

¹ Korea University
² Samsung Electronics

<https://ku-cvlab.github.io/Perturbed-Attention-Guidance>



Vezérelt Generálás

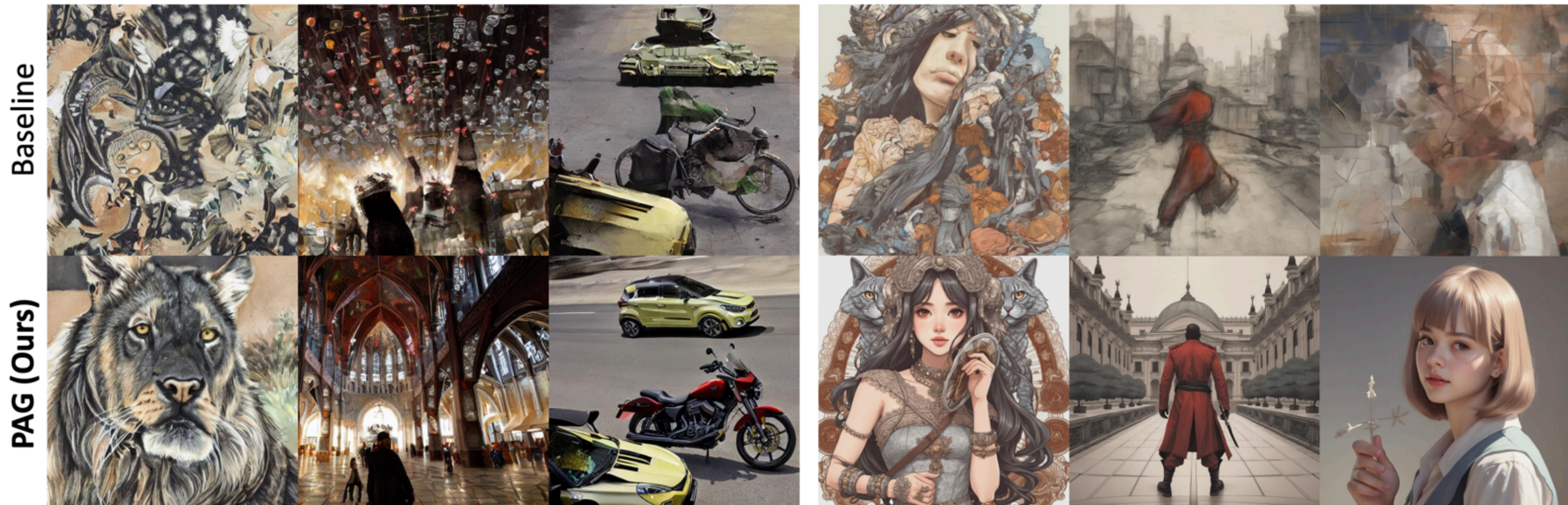
Perturbed-Attention Guidance (PAG)

Self-Rectifying Diffusion Sampling
with Perturbed-Attention Guidance

Donghoon Ahn⁺¹, Hyoungwon Cho⁺¹, Jaewon Min¹, Wooseok Jang¹,
Jungwoo Kim¹, SeonHwa Kim¹, Hyun Hee Park²,
Kyong Hwan Jin¹, and Seungryong Kim¹

¹ Korea University
² Samsung Electronics

<https://ku-cvlab.github.io/Perturbed-Attention-Guidance>



Stable Diffusion 1.5

SDXL

A PAG a feltétel nélkül generált képeken is sokat javít!

Vezérelt Generálás

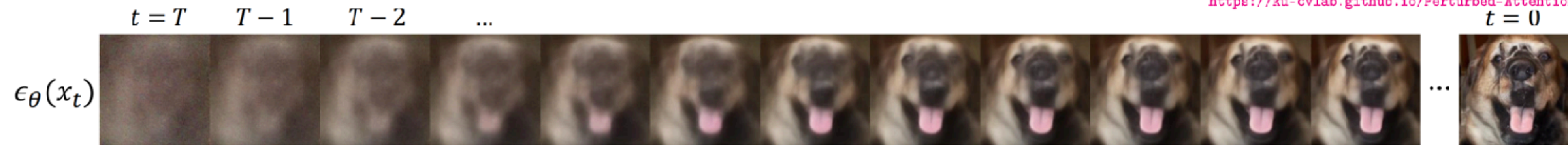
Perturbed-Attention Guidance (PAG)

Self-Rectifying Diffusion Sampling
with Perturbed-Attention Guidance

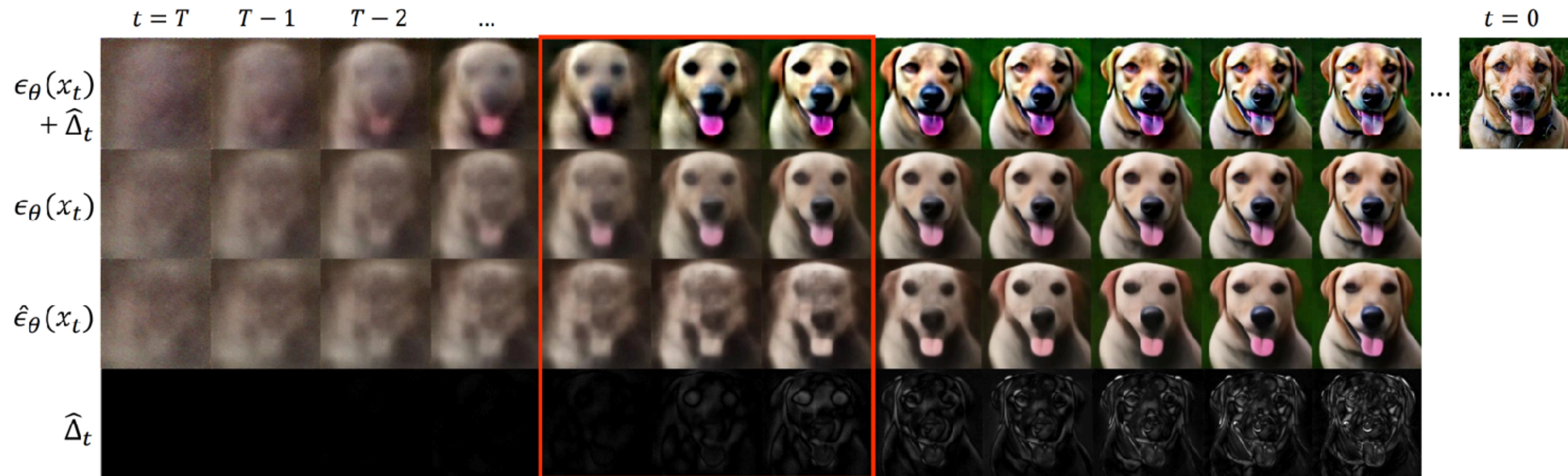
Donghoon Ahn^{*1}, Hyoungwon Cho^{*1}, Jaewon Min¹, Wooseok Jang¹,
Jungwoo Kim¹, SeonHwa Kim¹, Hyun Hee Park²,
Kyong Hwan Jin¹, and Seungryong Kim¹

¹ Korea University
² Samsung Electronics

<https://ku-cvlab.github.io/Perturbed-Attention-Guidance>



(a) Sampling process without PAG



(b) Sampling process with PAG

Vezérelt Generálás

Promptolás képpel

Képi promptok



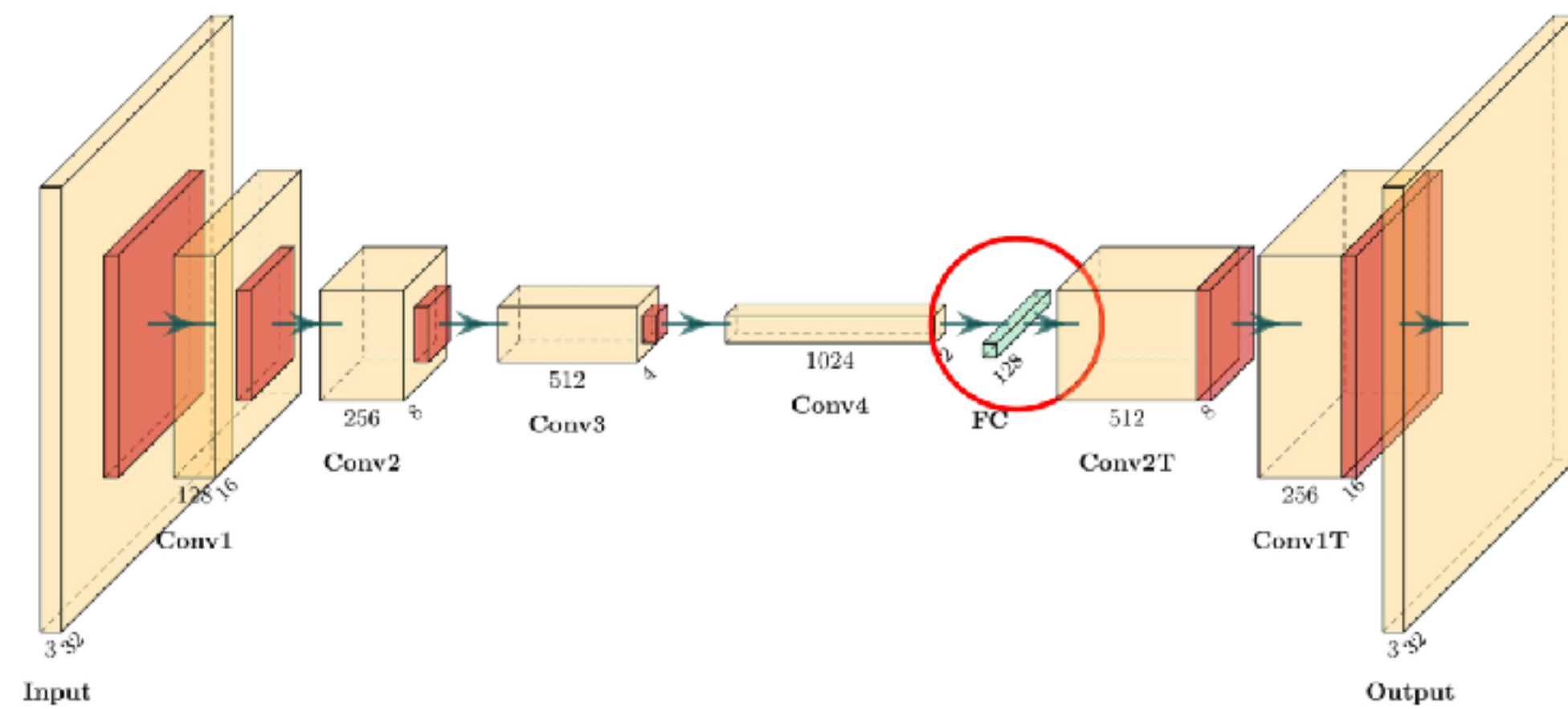
Hogyan
kódoljunk be
képi feltételeket?

Generált képek

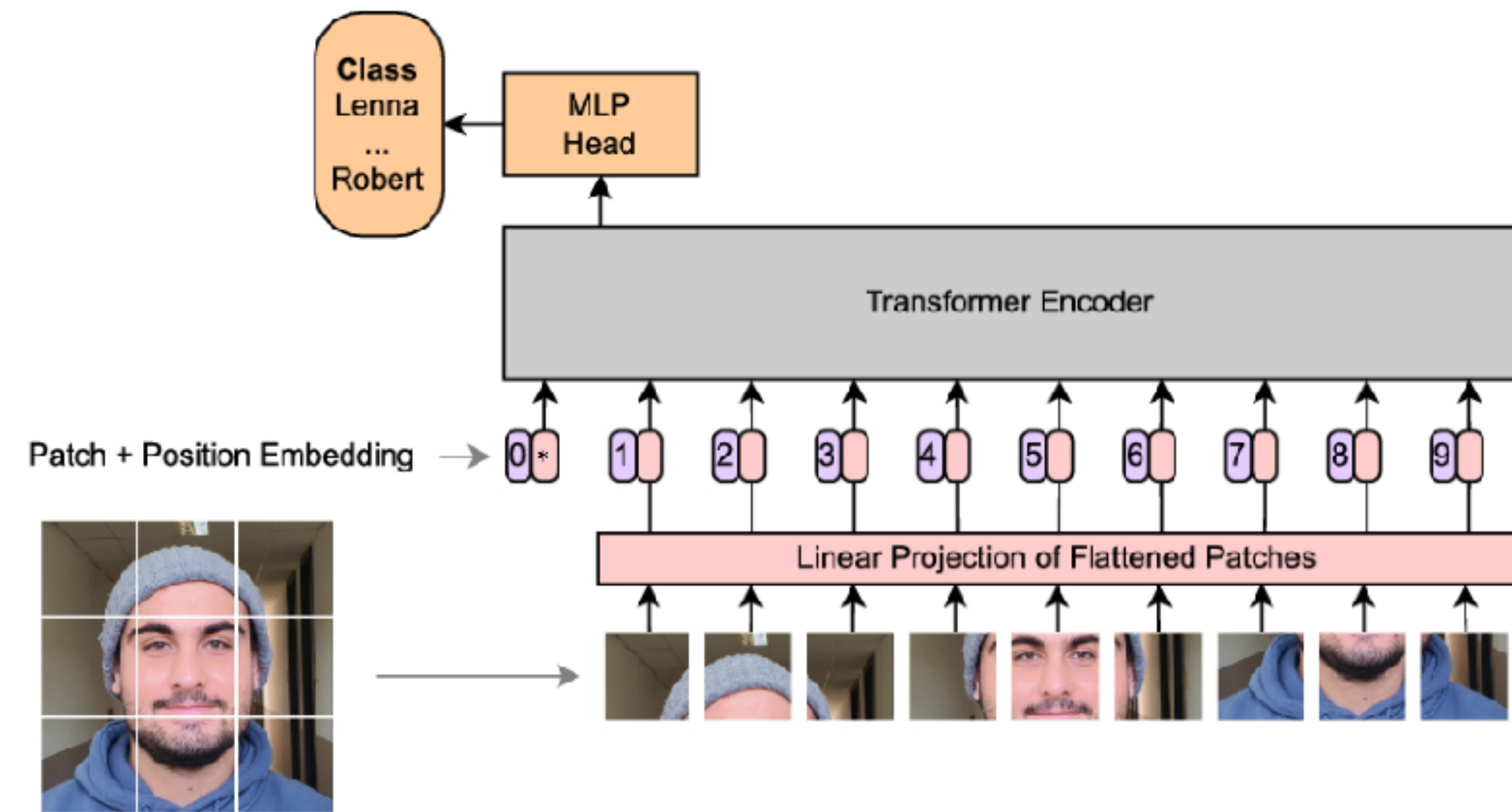


Képek Kódolása

VAE, ViT



Konvolúció (V)AE



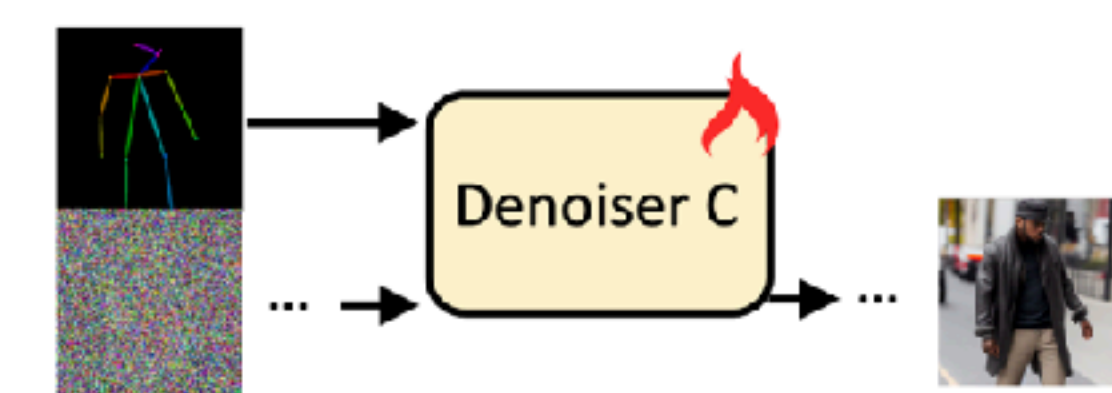
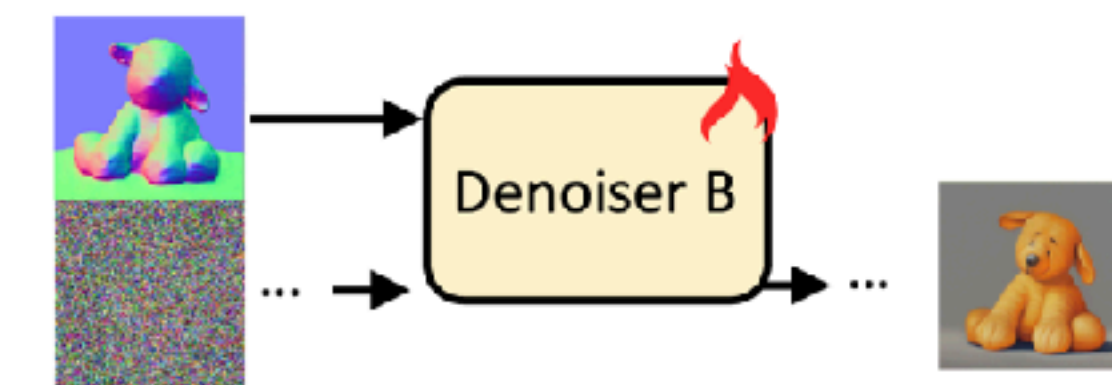
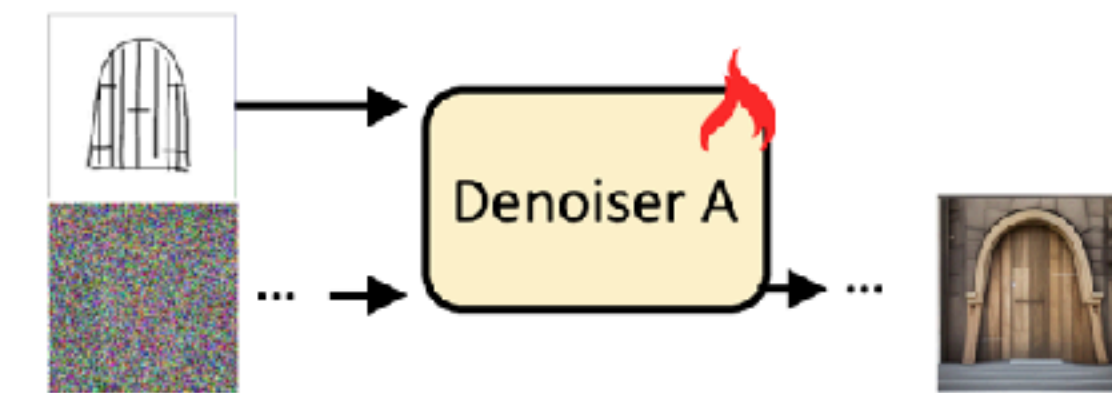
Vision Transformer (ViT)

(pl. CLIP / MAE / DINO)

Finomhangolás

Képi kondíciók figyelembe vétele

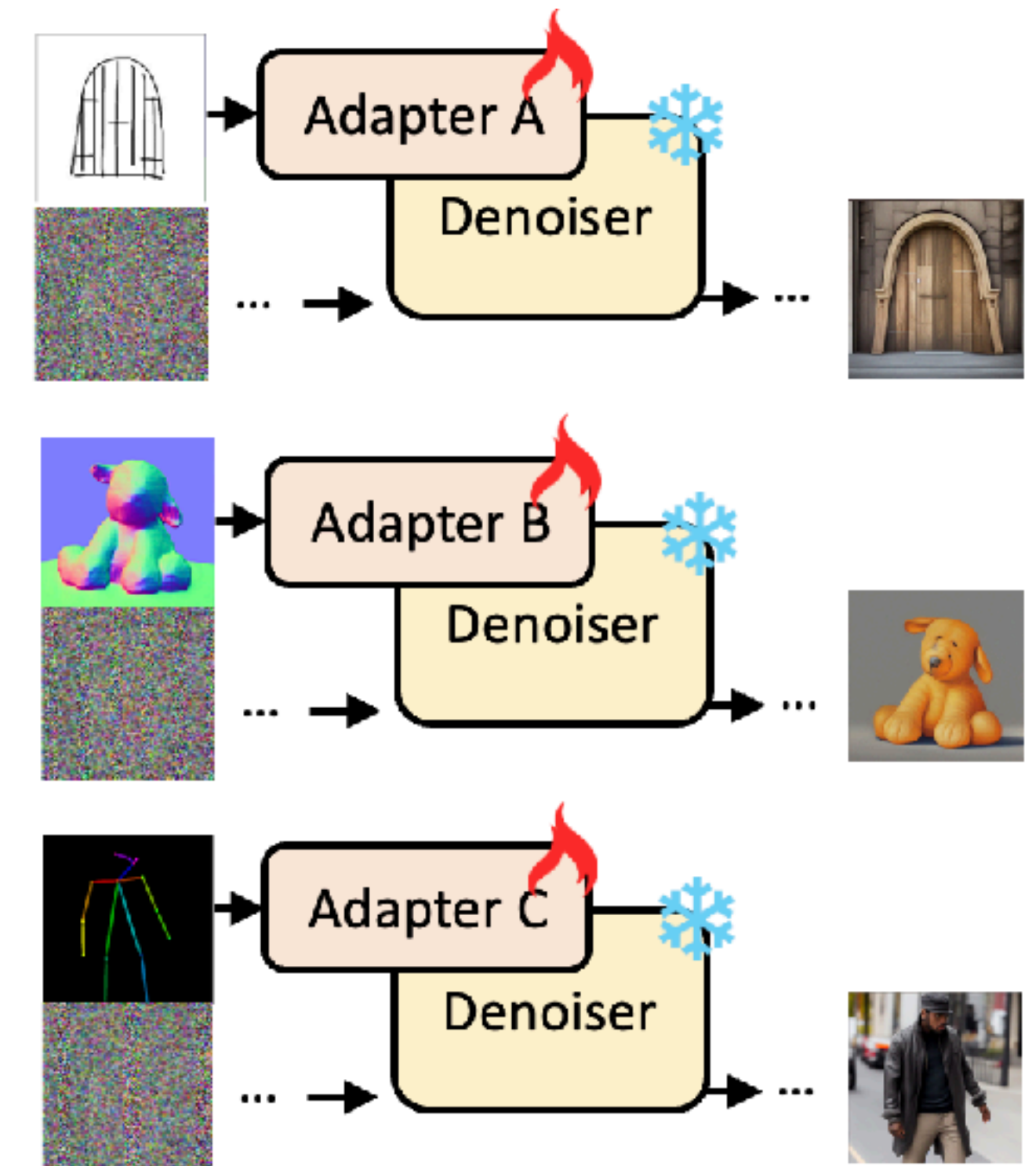
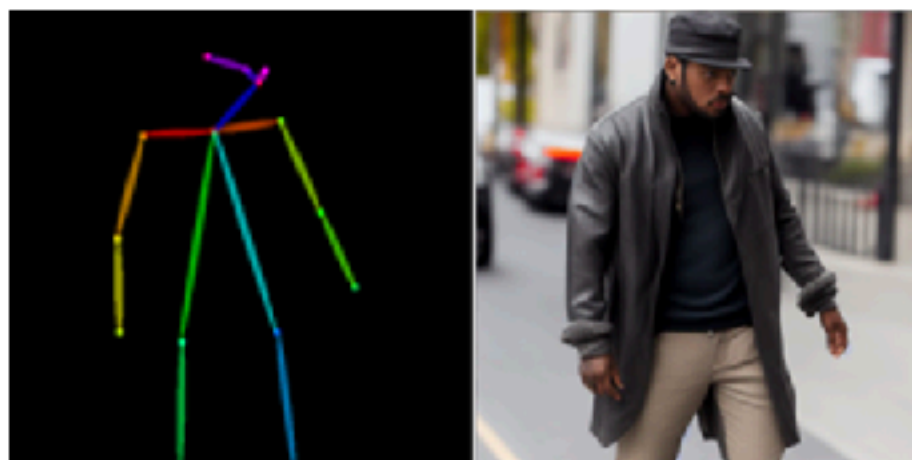
- Hogyan tanítsunk meg egy diffúziós modellt korábban nem látott kondíciókra — pl. képekre?
- Adjuk hozzá az új kondíciót a korábbiakhoz, aztán **finomhangoljuk** példák alapján?
 - Rengeteg paramétert kéne hangolni (pl. SD 3.5: 8B!)
 - Kevés lesz az adat (az eredeti hálót akár több *milliárd* képen tanították...)
 - Minden újabb kondíció komoly idő és compute ráfordítást igényelne
 - Összezavarhatjuk vele az eredeti hálót (katasztrofális felejtés)



Finomhangolás

Adapterek

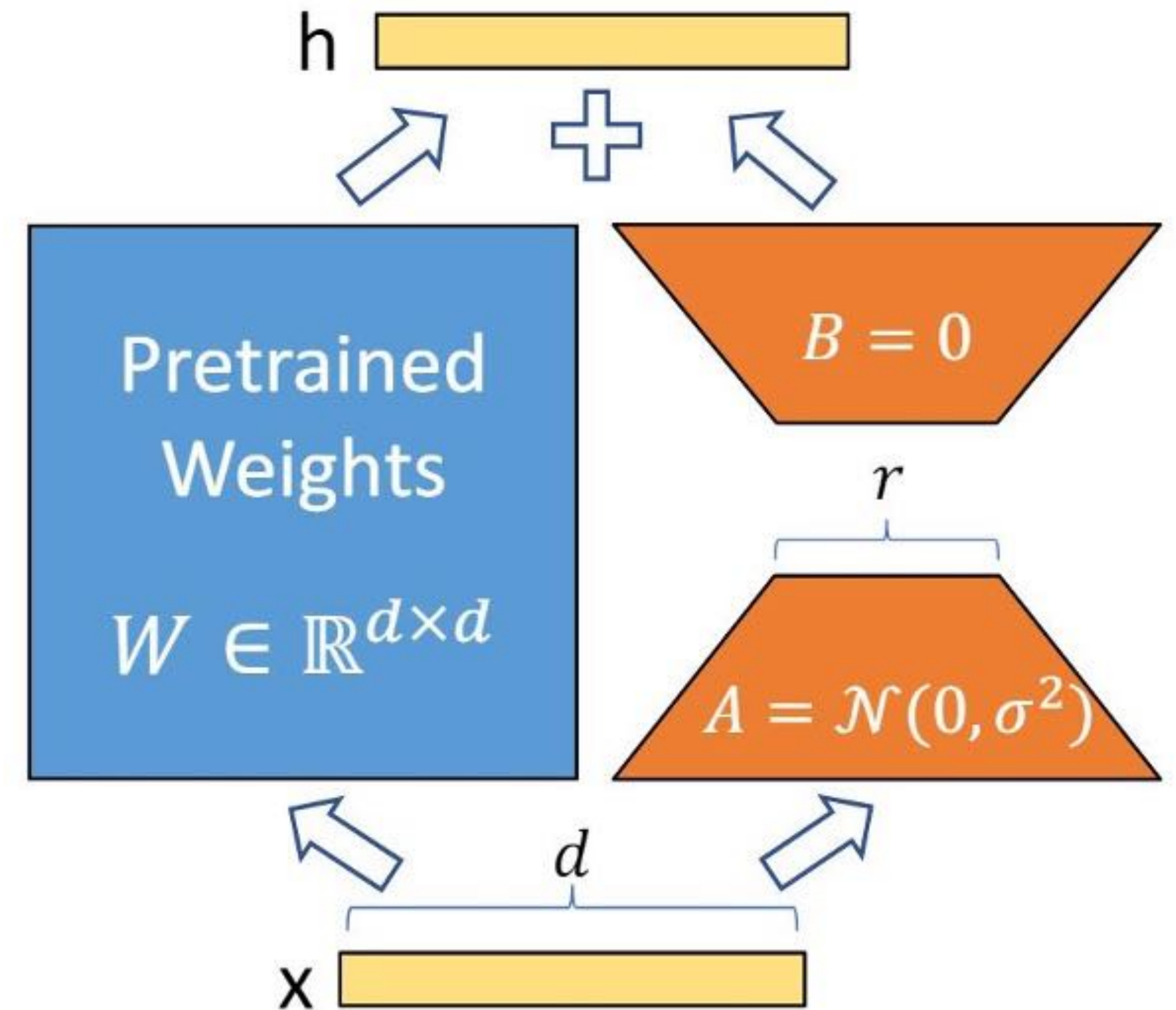
- Ötlet: az eredeti nagy diffúziós hálót “fagyasszuk be”, csak egy (jóval kisebb) adapter háló paramétereit tanítsunk!
- Népszerű adapterek:
 - LoRA
 - ControlNet
 - IP-Adapter



Finomhangolás

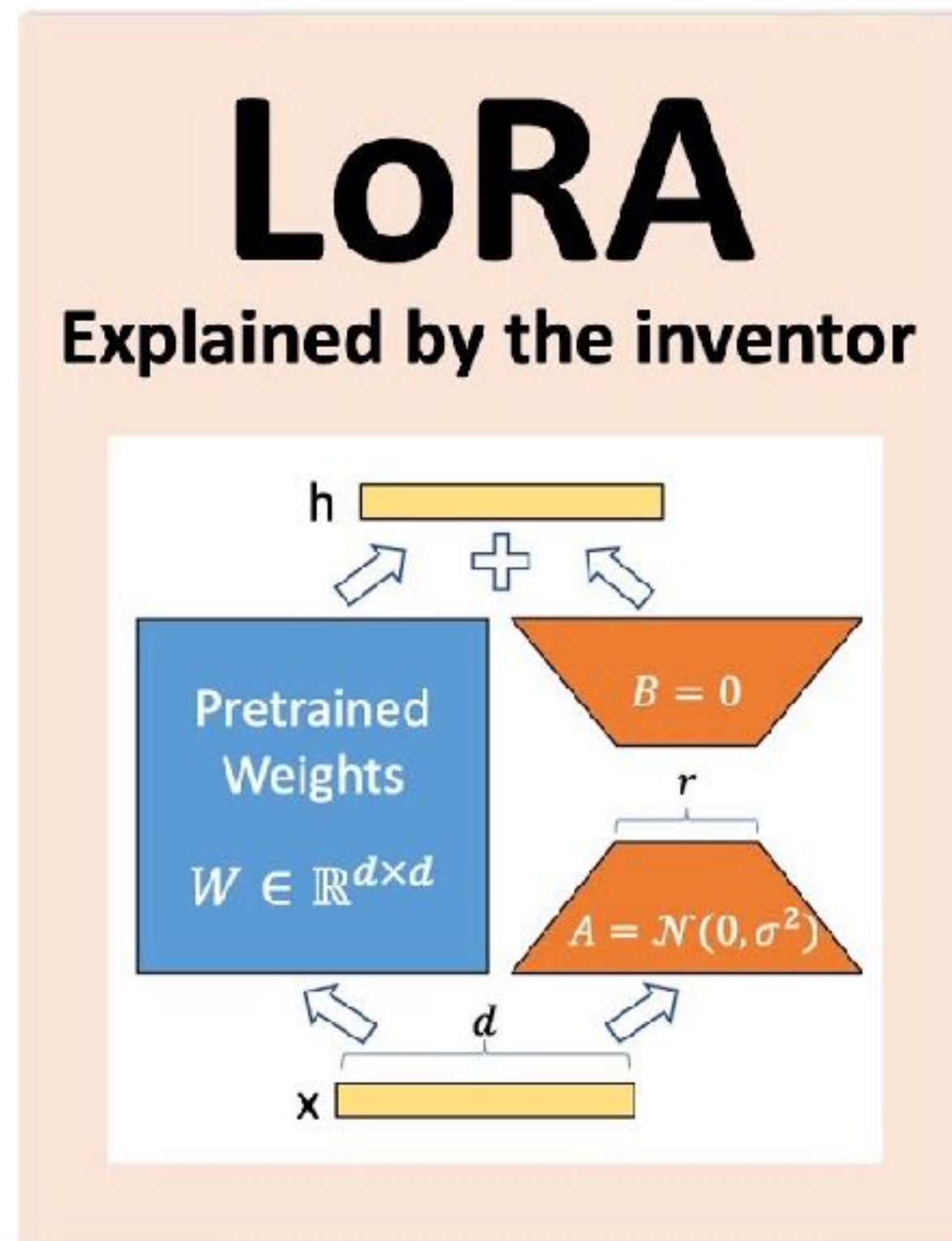
Adapterek – LoRA

- Megfigyelés: finomhangolás során a hálóparaméterek (pl. QKV mátrixok) alacsony dimenziós altérben változnak!
- Paraméterezzük a módosításokat alacsony rangú mátrixként (külső szorzatként) – **Low-Rank Adaptation (LoRA)**
- A LoRA mátrixok kis méretűek, gyorsan taníthatóak és több LoRA akár *kombinálható is!*
- Nagyon népszerű módszer – főleg LLM-ek és képgenerátorok attention rétegeinek hangolására
 - Nagy modellek finomhangolása (pl. A többi adapter implementációja is) általában LoRA-val történik!



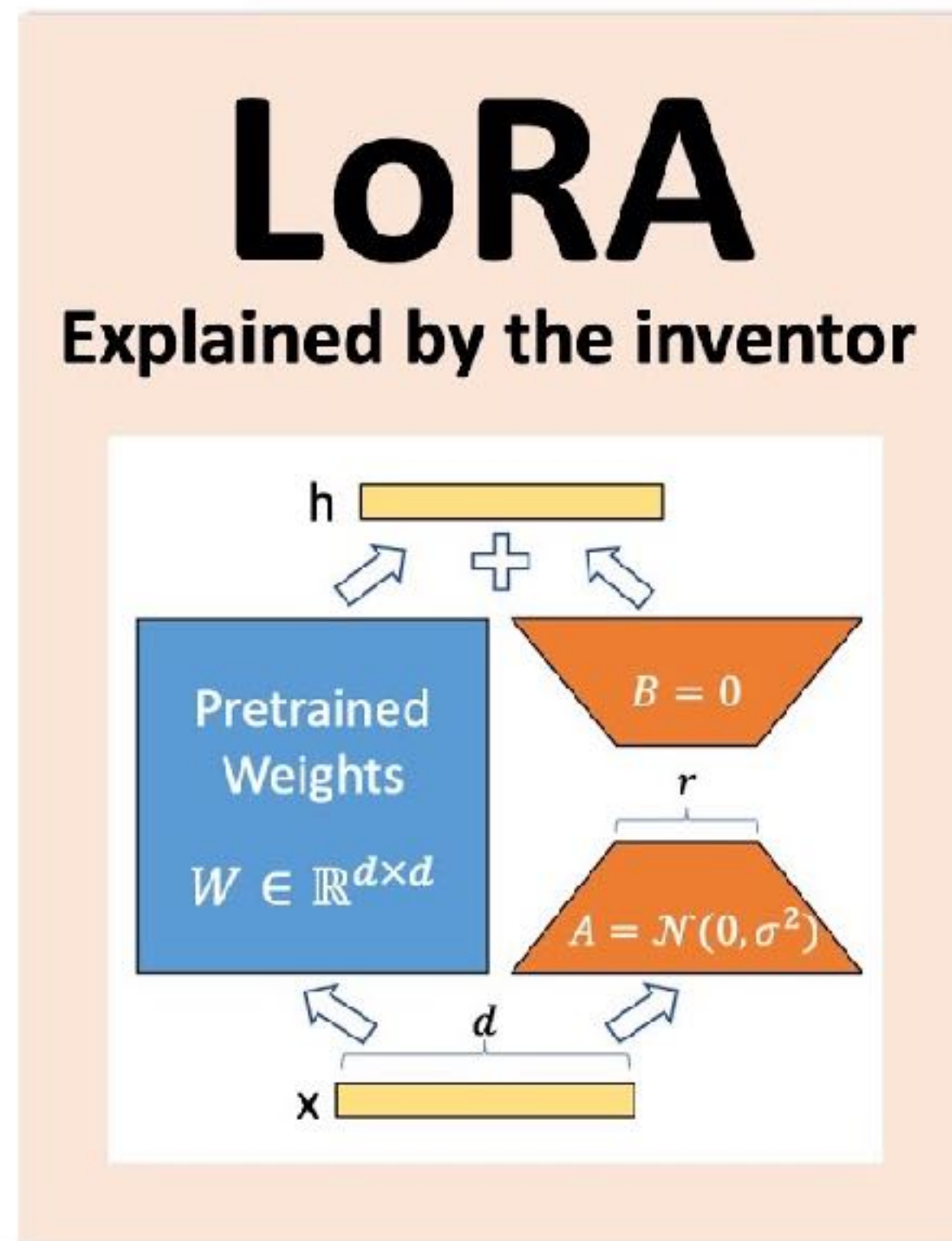
Finomhangolás

Adapterek – LoRA – Videóajánló



Finomhangolás

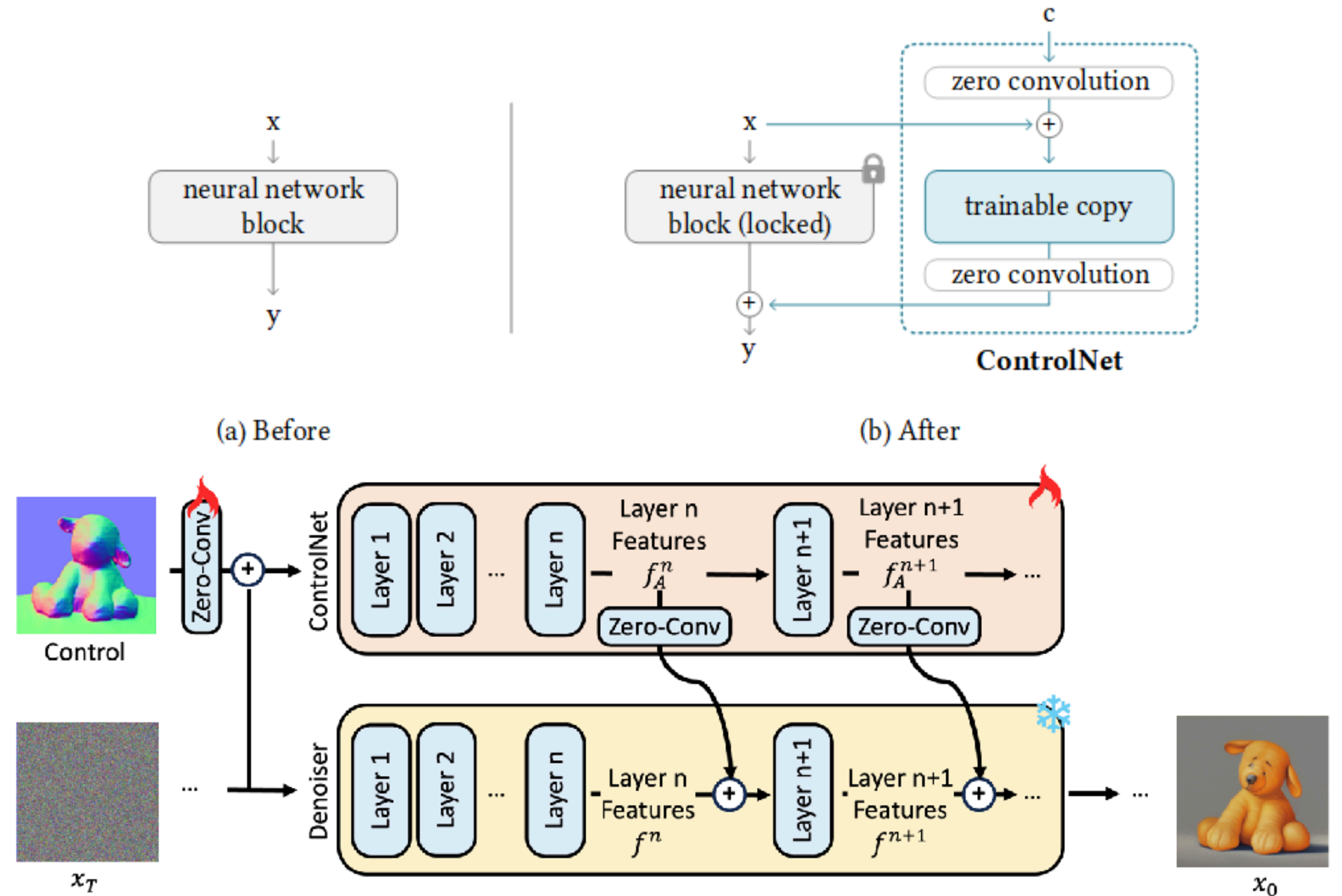
Adapterek – LoRA – Videóajánló



Finomhangolás

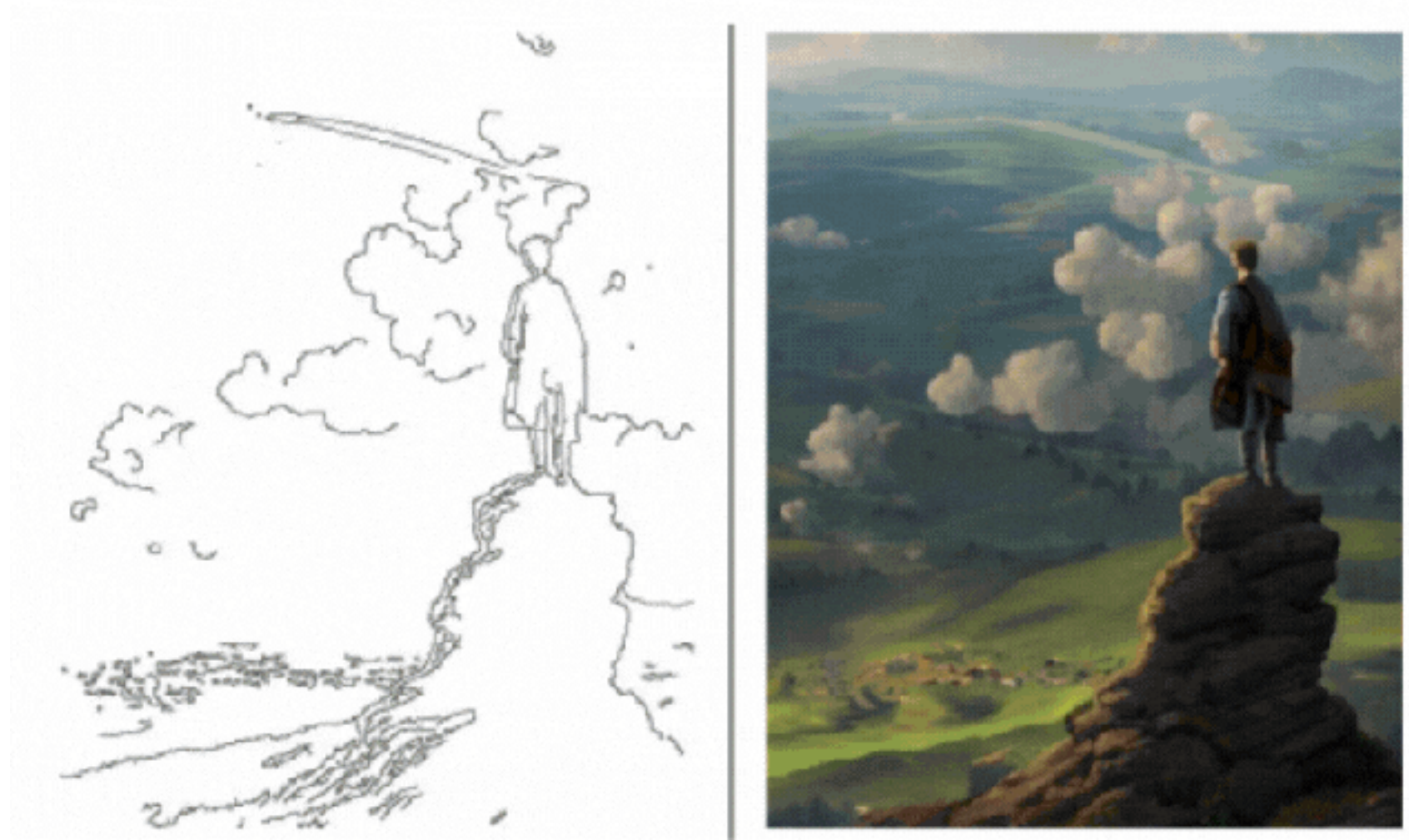
Adapterek – ControlNet

- Finomhangoljuk a háló másolatát, és adjuk hozzá az eredetihez – **ControlNet**
- Kezdetben a másolatnak ne legyen hatása – nullára inicializált konvolúciós (Zero-Conv) rétegekkel adjuk össze!



Finomhangolás

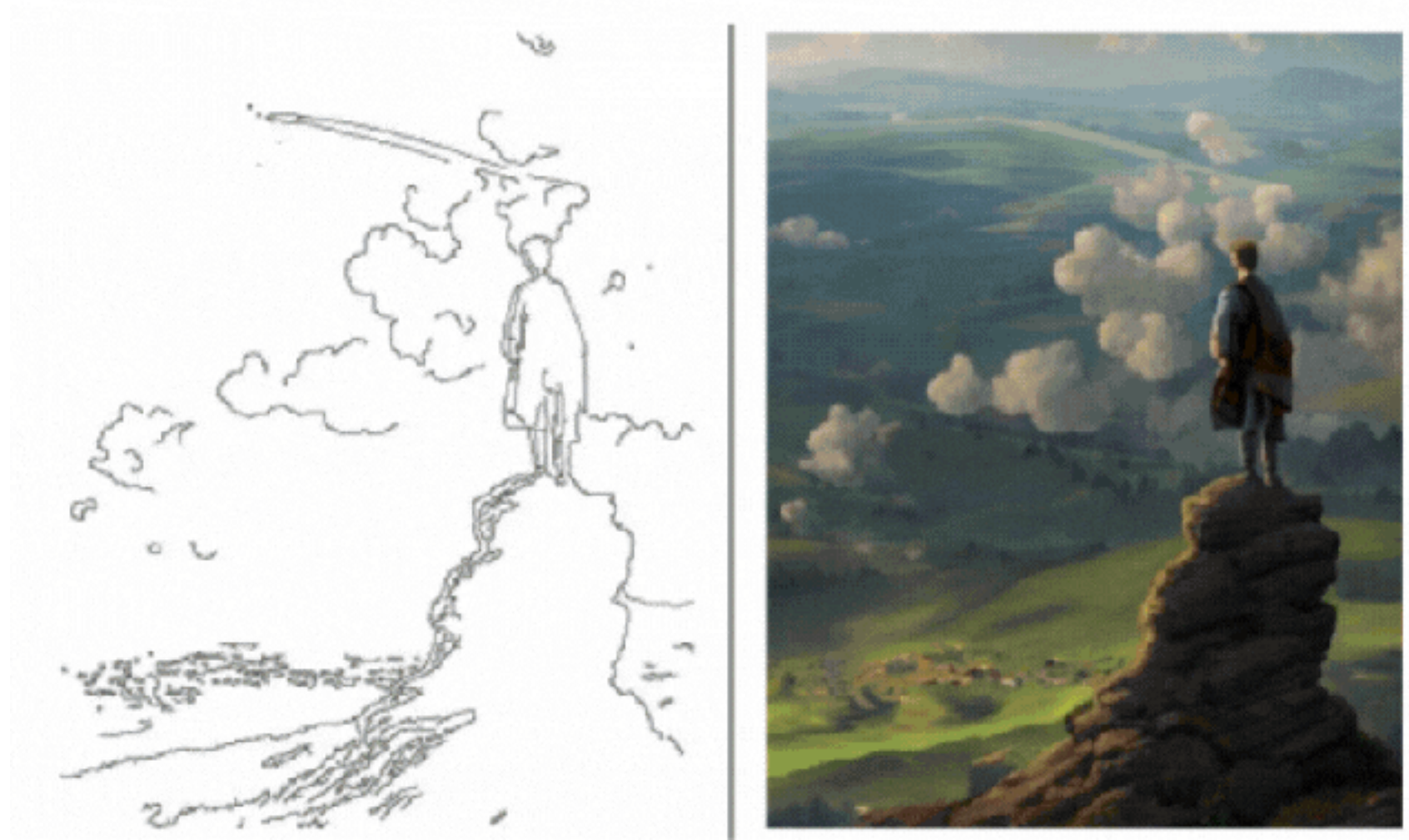
Adapterek – ControlNet



Sketch	Normal map	Depth map	Canny[11] edge	M-LSD[24] line	HED[91] edge	ADE20k[96] seg.	Human pose

Finomhangolás

Adapterek – ControlNet



Sketch	Normal map	Depth map	Canny[11] edge	M-LSD[24] line	HED[91] edge	ADE20k[96] seg.	Human pose

Finomhangolás

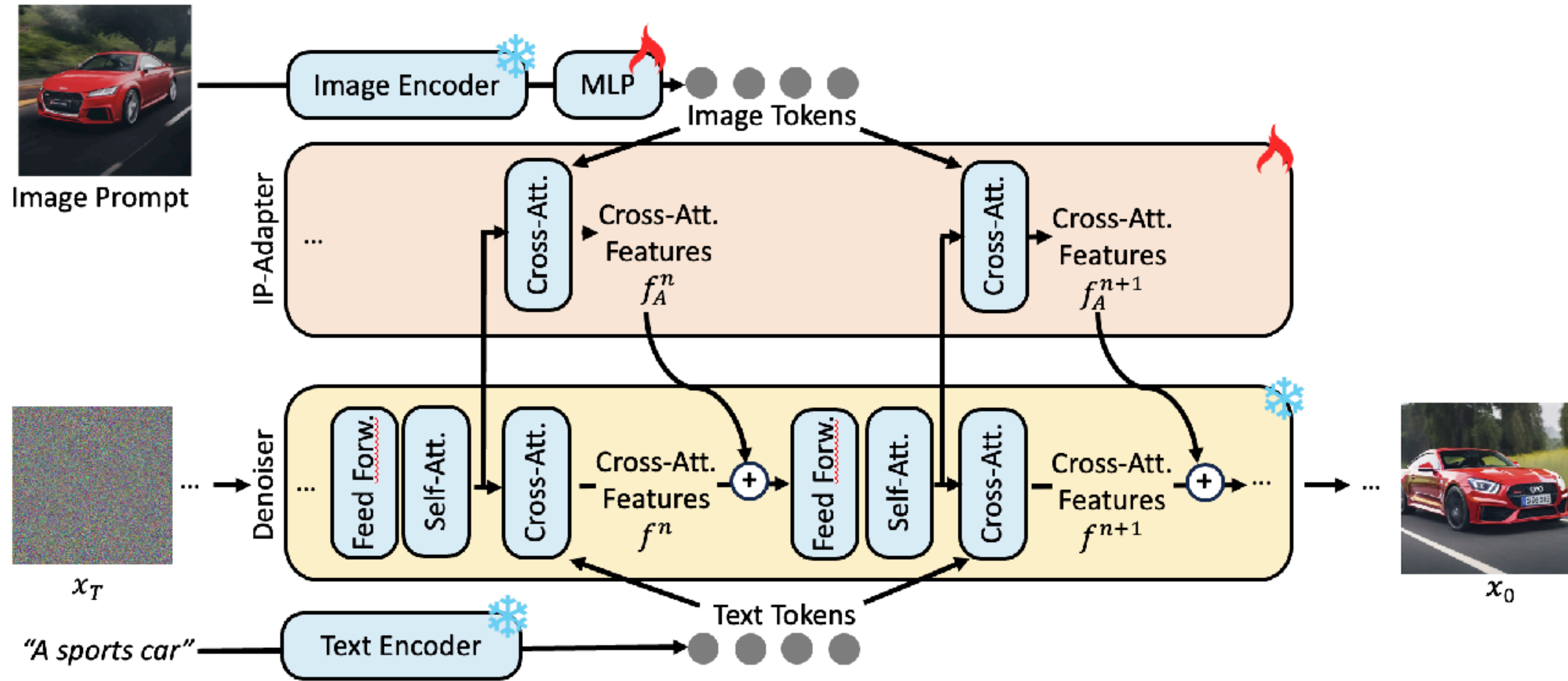
Adapterek – ControlNet



ControlNet (szöveges prompt nélkül) + PAG

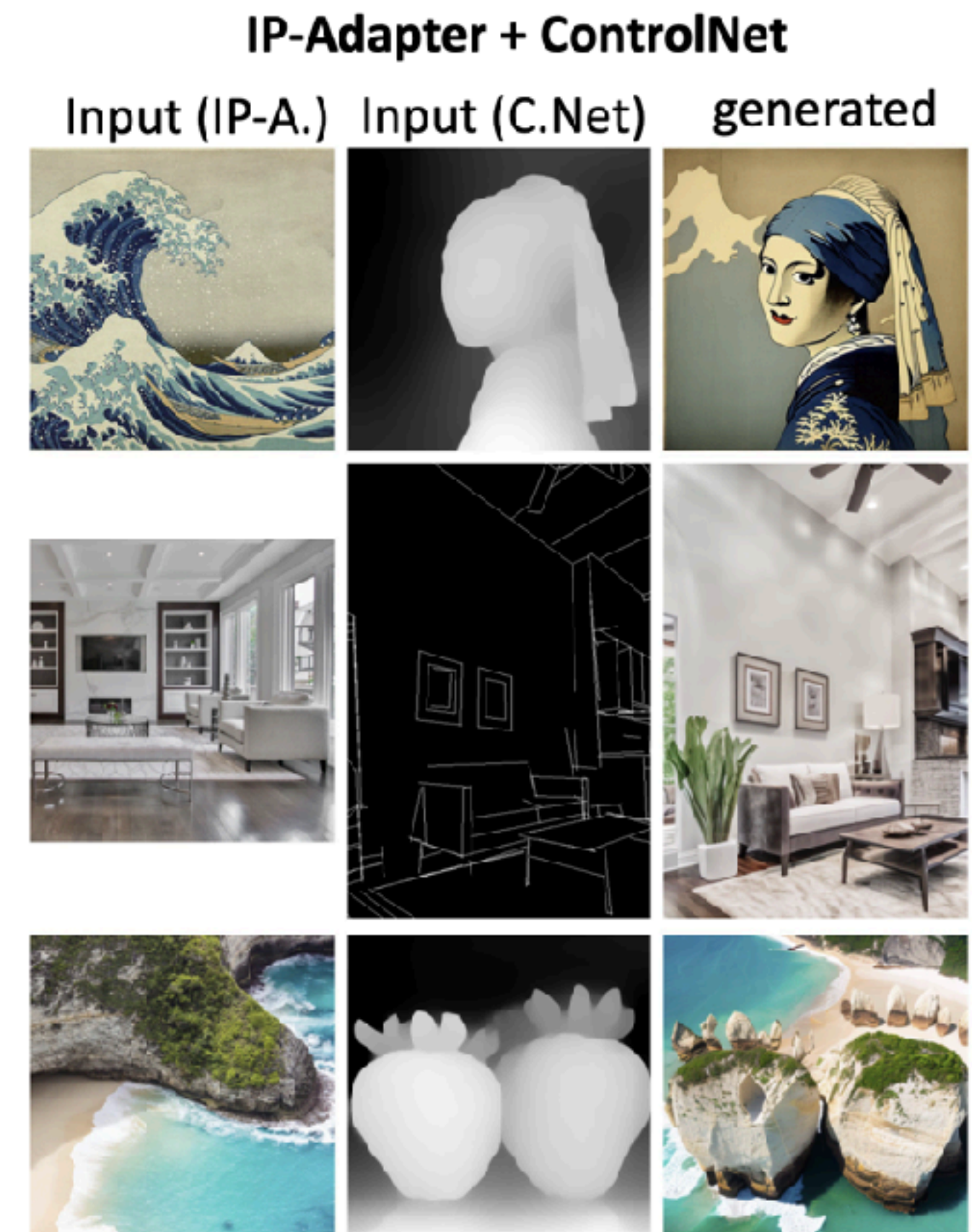
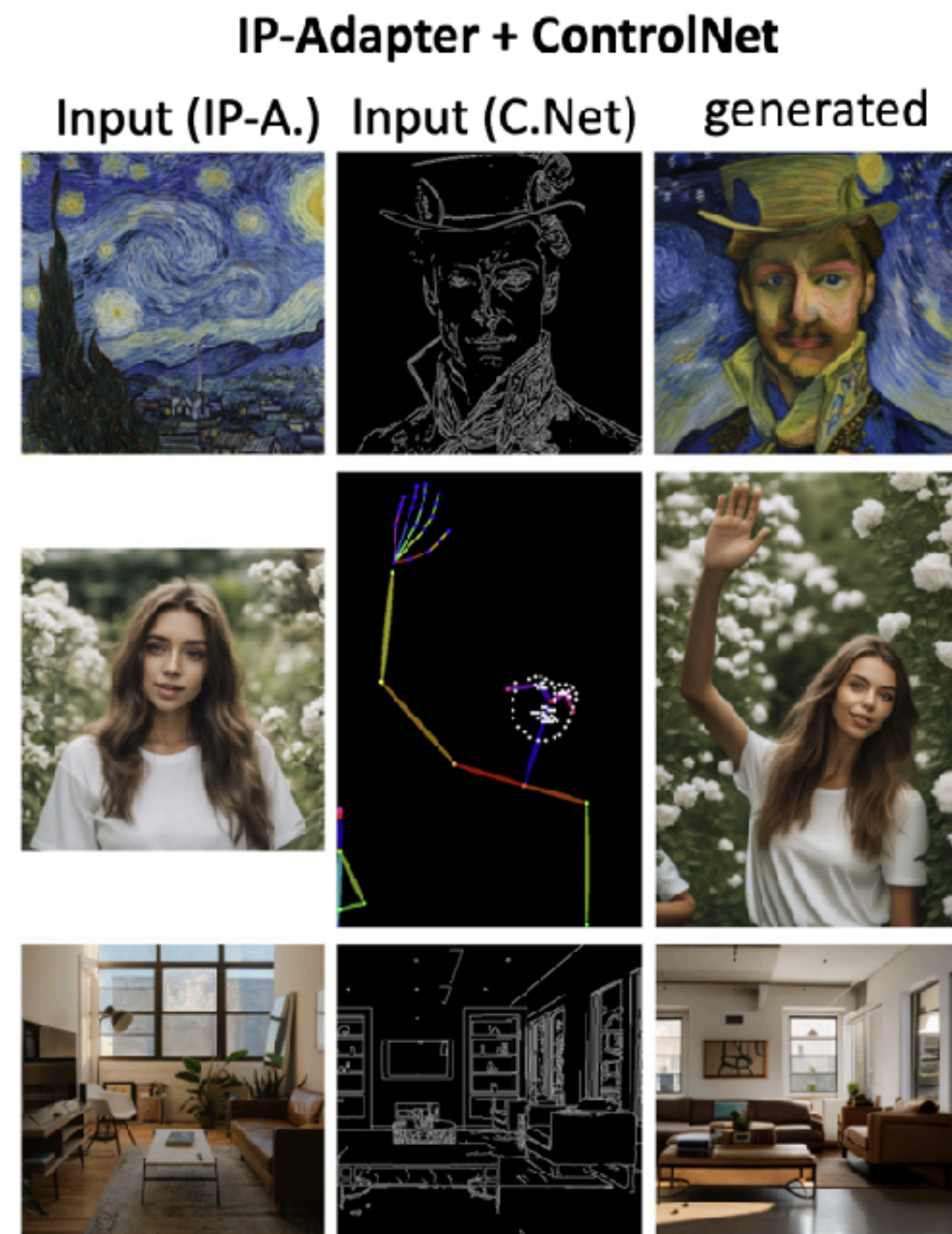
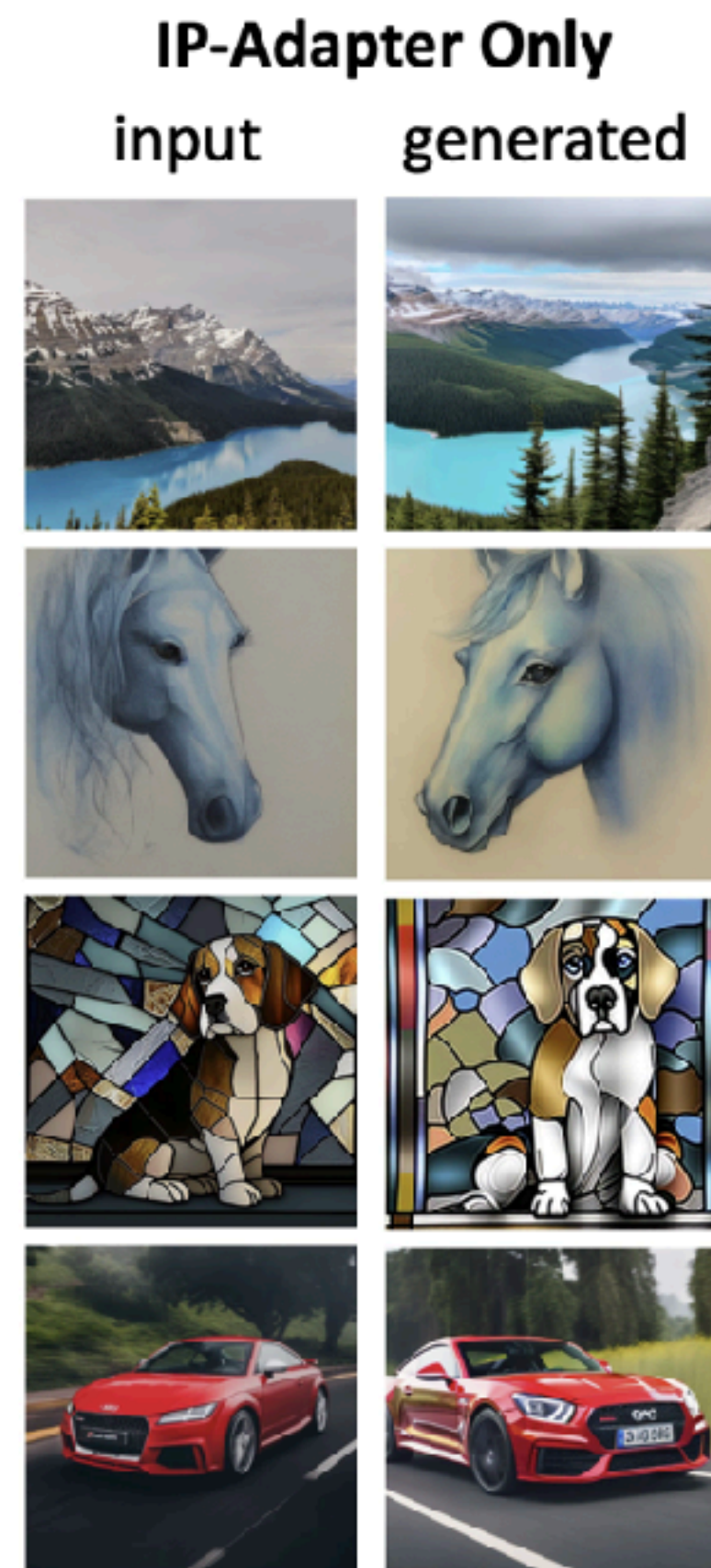
Finomhangolás

Adapterek – IP-Adapter



Finomhangolás

Adapterek – IP-Adapter



Finomhangolás

DreamBooth

DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation

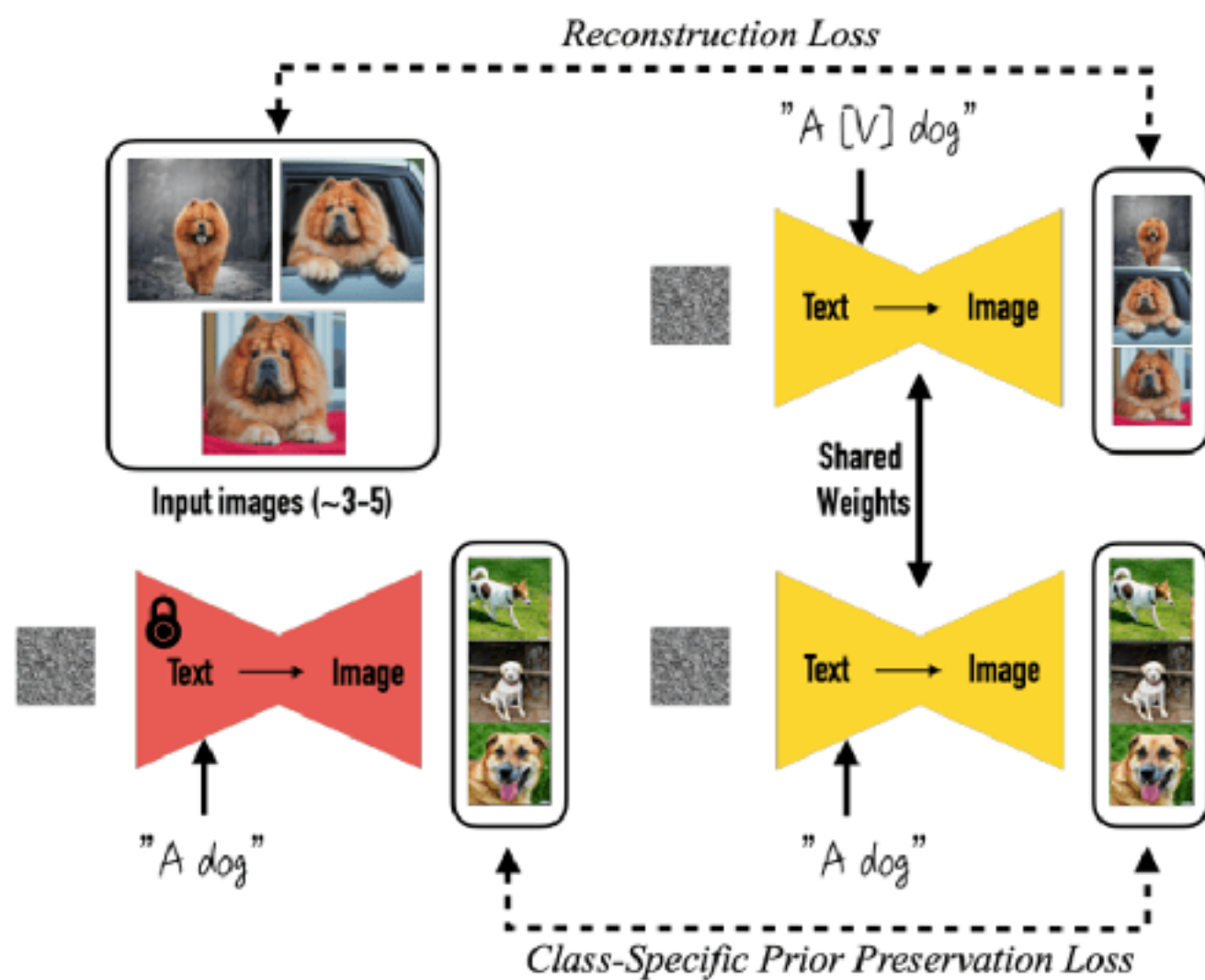
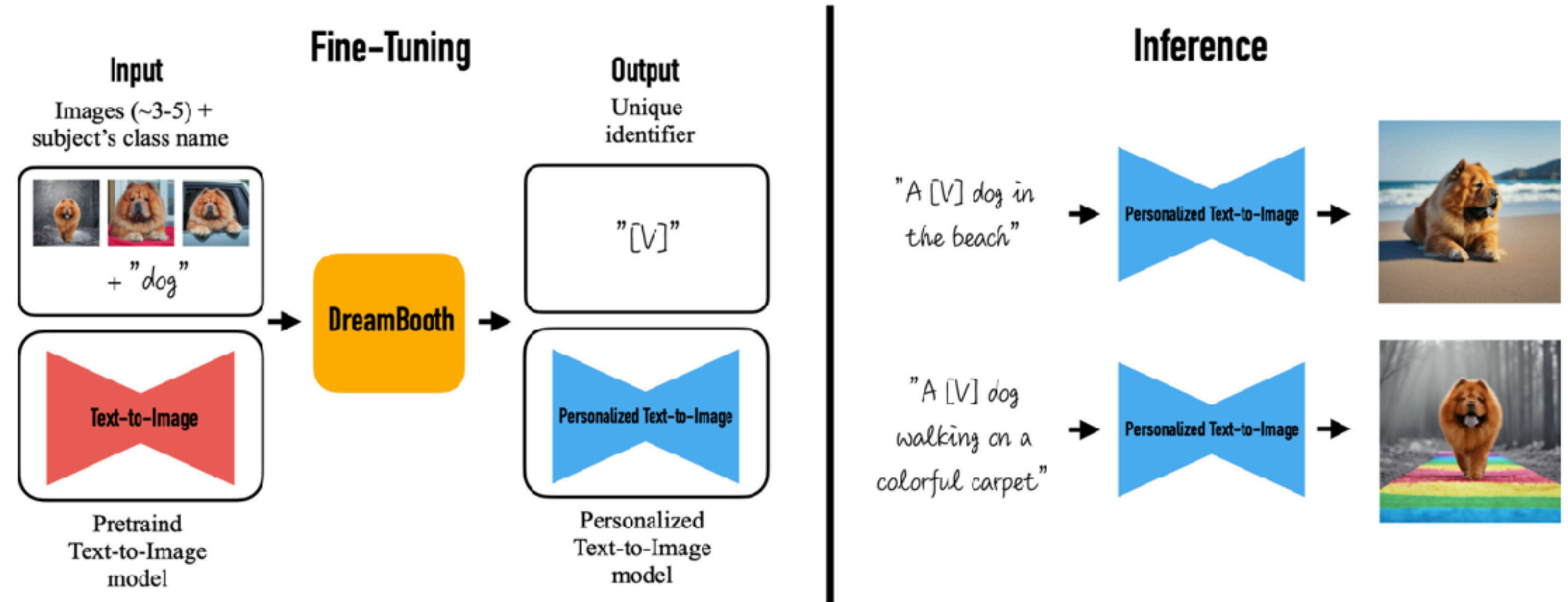
Nataniel Ruiz^{*,1,2}
Yael Pritch¹

Yuanzhen Li¹
Michael Rubinstein¹

Varun Jampani¹
Kfir Aberman¹

¹ Google Research ² Boston University

- **DreamBooth**: finomhangolás kevés (3-5) példakép alapján!
- Eredmény: finomhangolt modell + specializált szöveges token



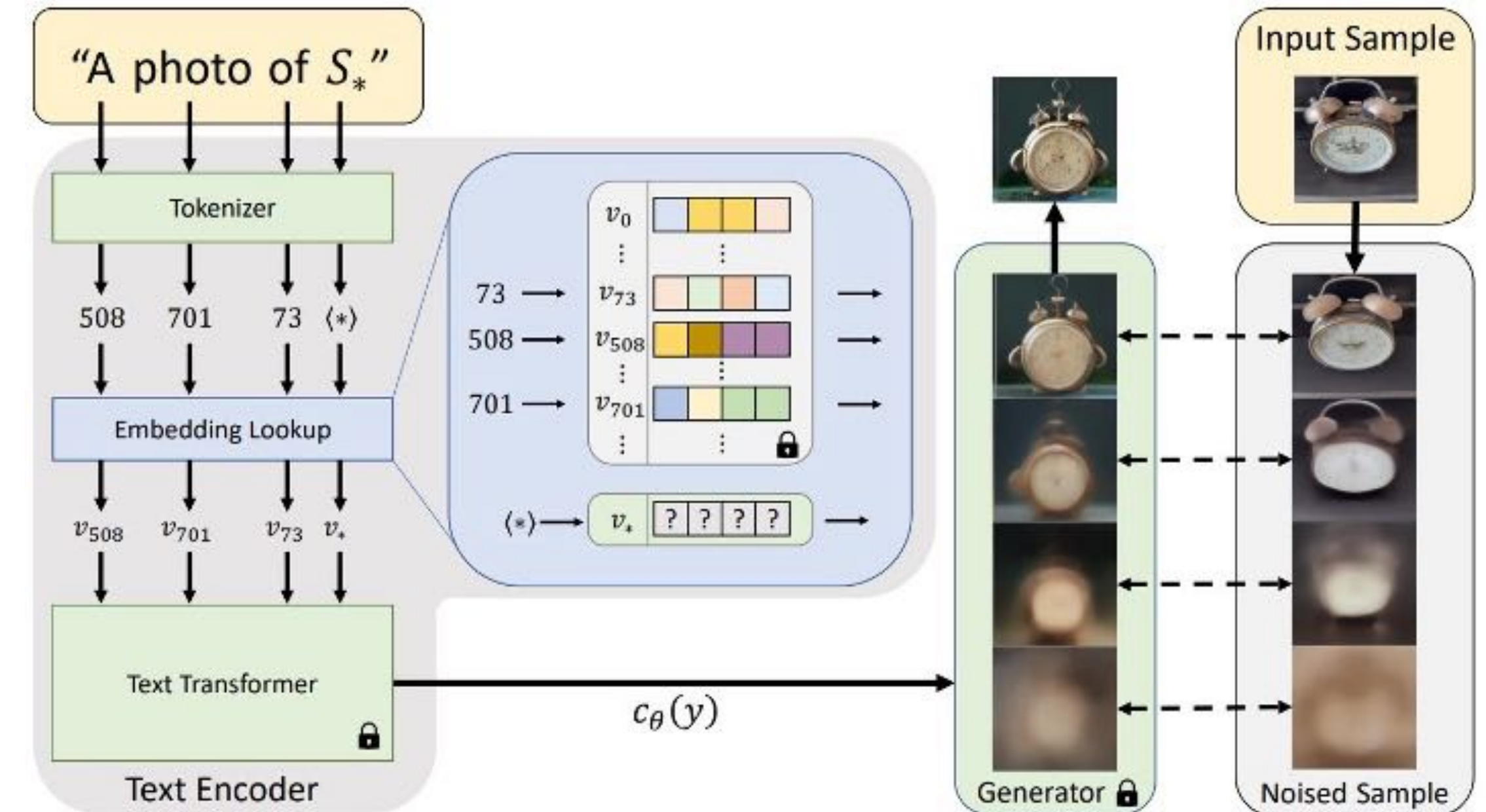
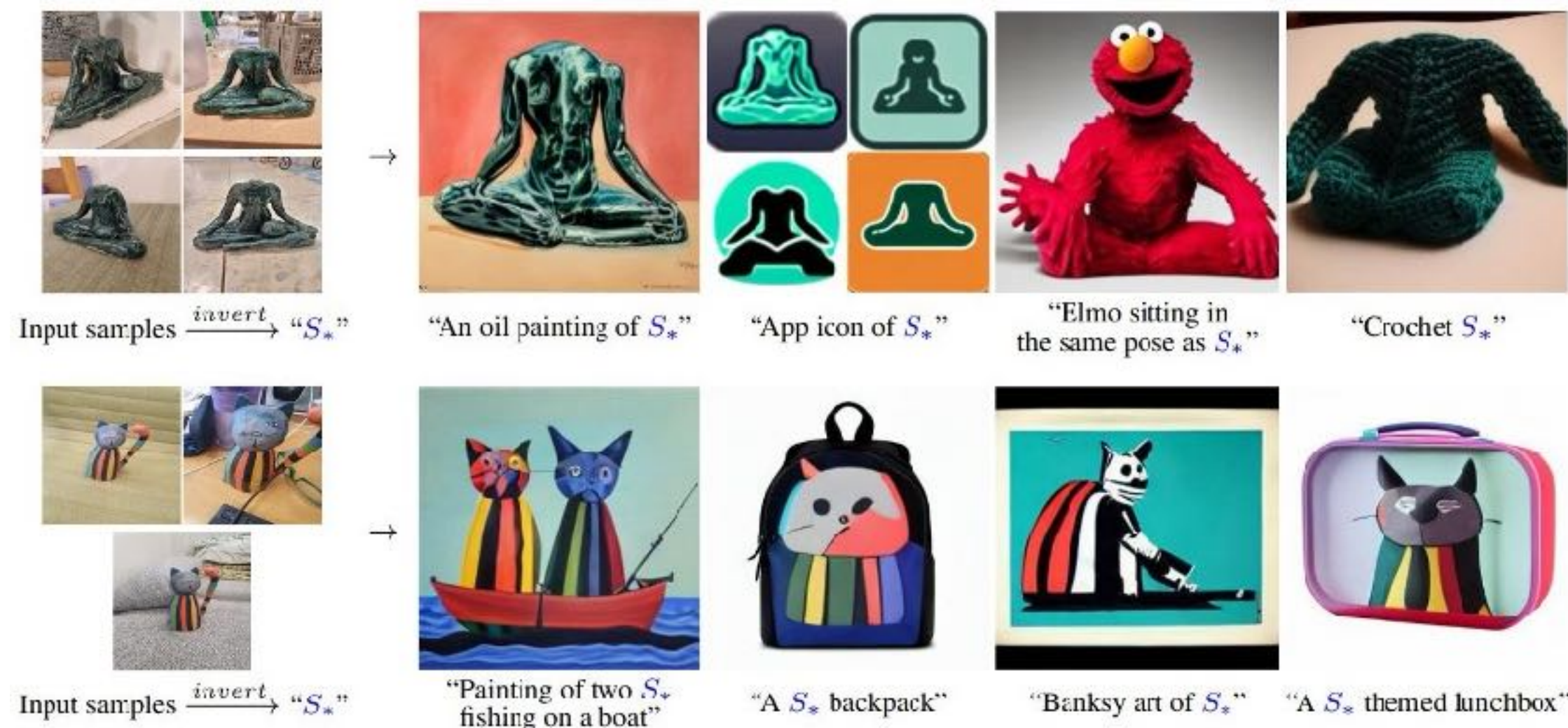
Finomhangolás

Szöveges inverzió

- Szöveges Inverzió – a példaképek optimalizálunk egy egyedi szöveges token!
- A képgenerátort **nem** finomhangoljuk, *csak a token*!

An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion

Rinon Gal^{1,2}, Yuval Alaluf¹, Yuval Atzmon², Or Patashnik¹, Amit H. Bermano¹, Gal Chechik²,
Daniel Cohen-Or¹,
¹Tel Aviv University, ²NVIDIA



Finomhangolás

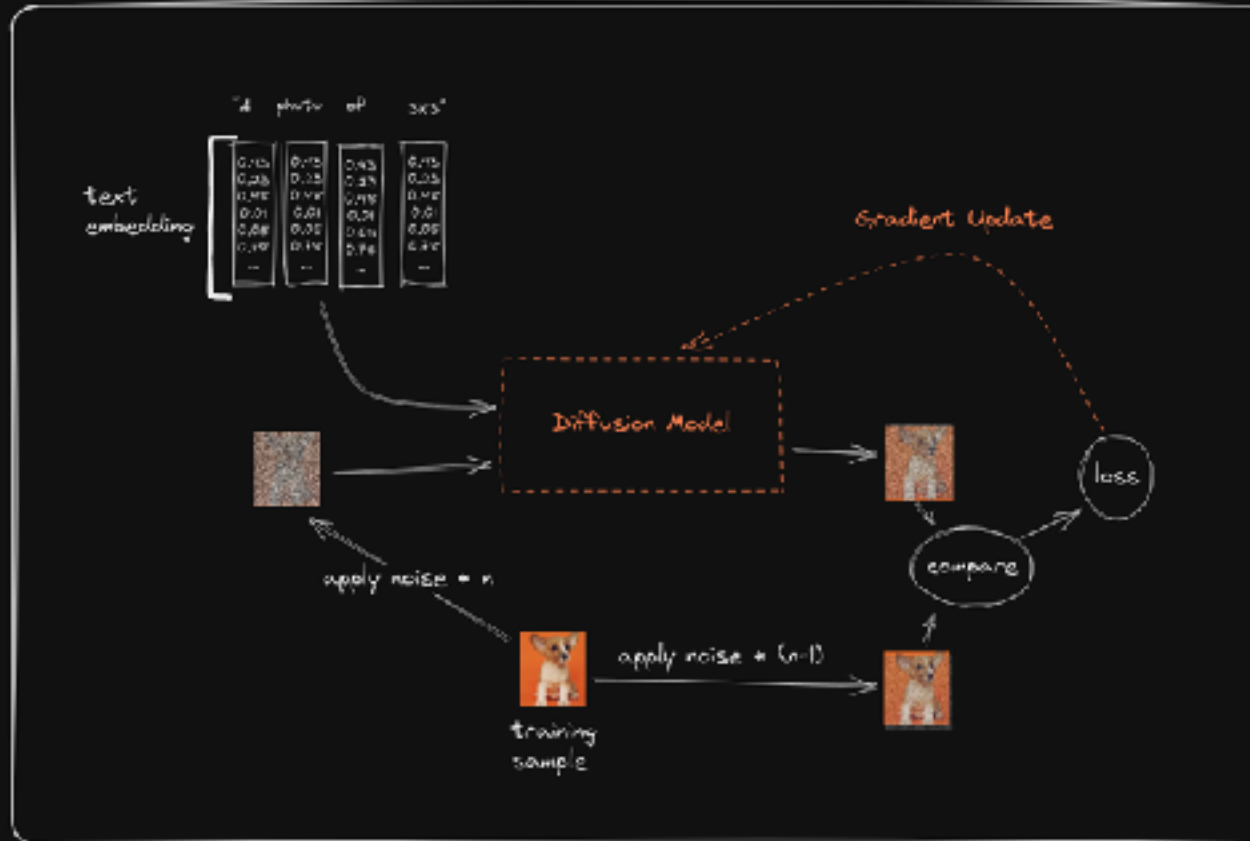
Áttekintés



Dreambooth

Fine-tunes the diffusion model itself until it understands the new concept

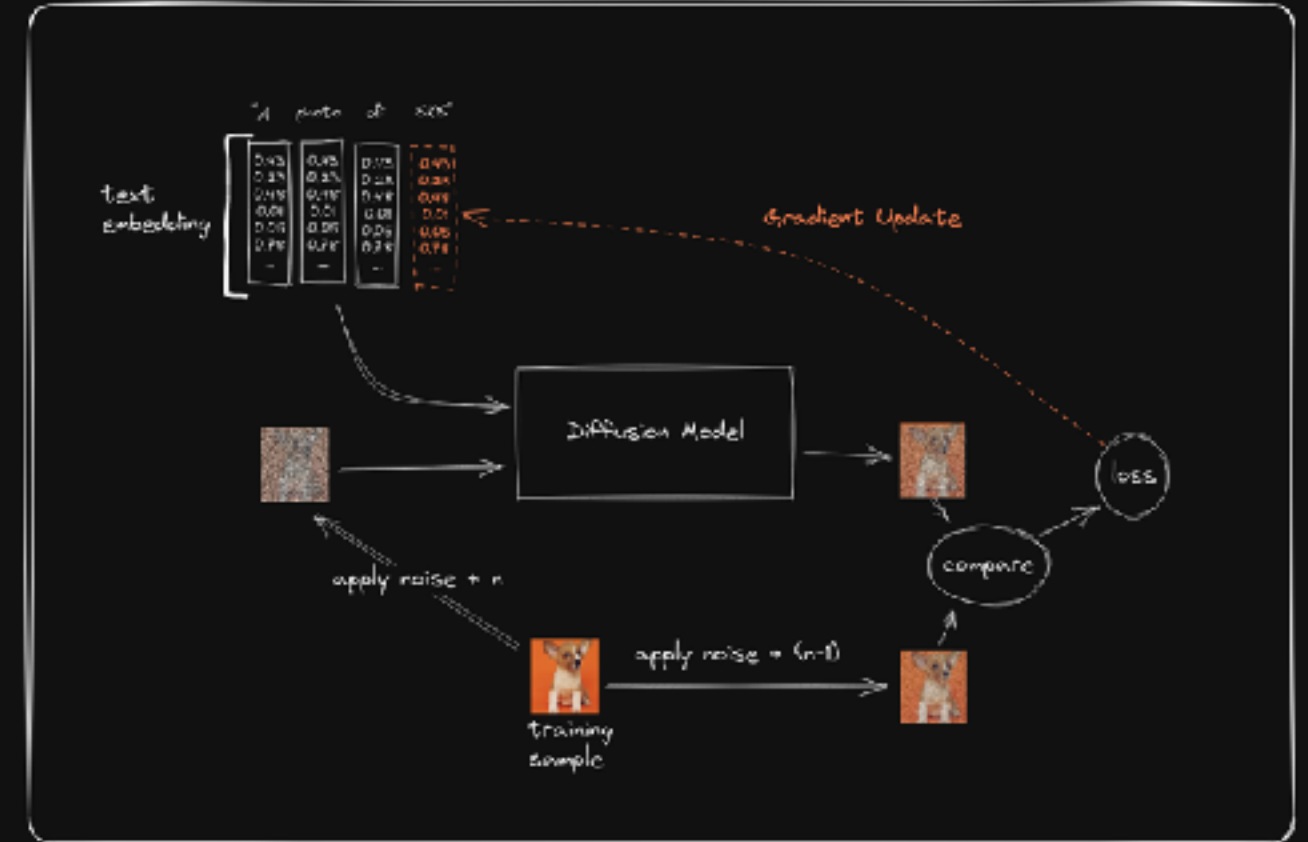
+ probably the most effective
- storage inefficient (whole new model to deal with)



Textual Inversion

Creates a special word embedding which captures the new concept

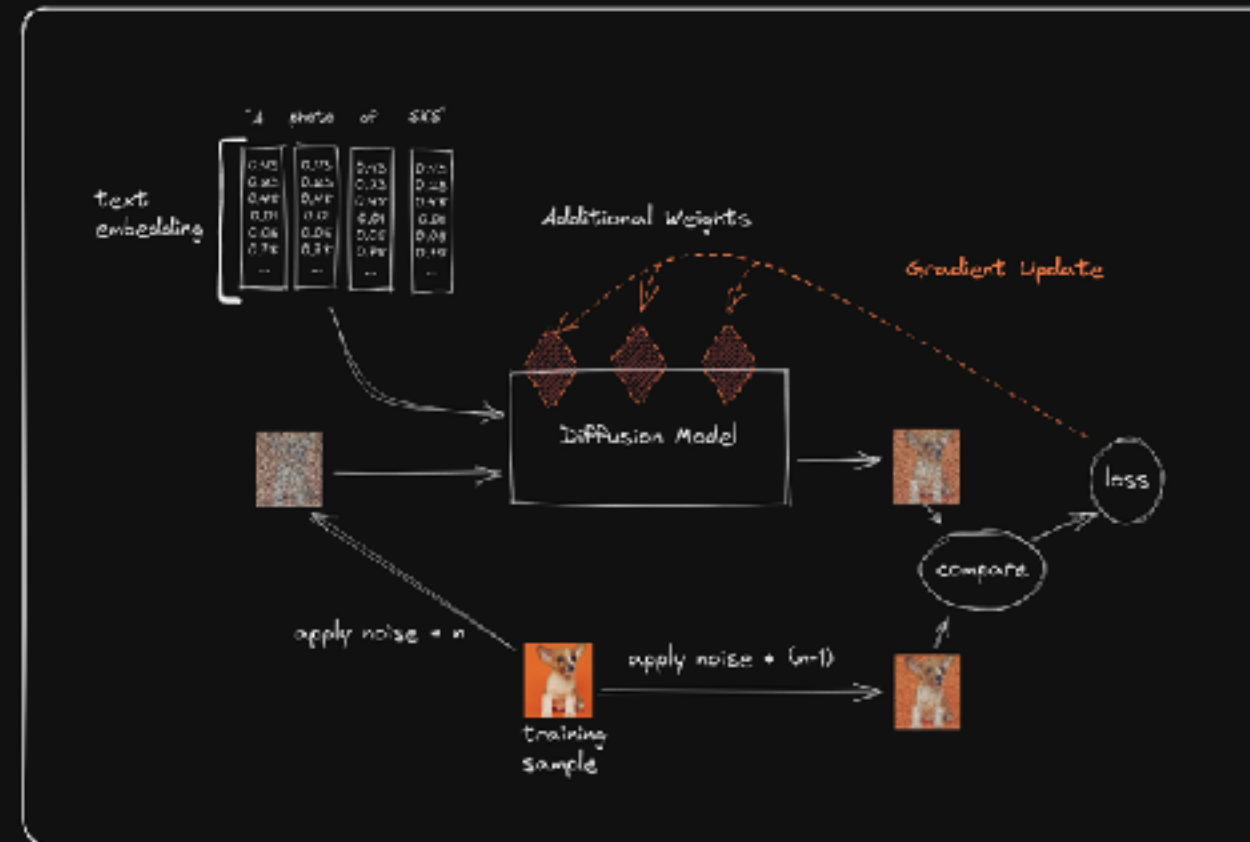
+ output is a tiny embedding



LoRA

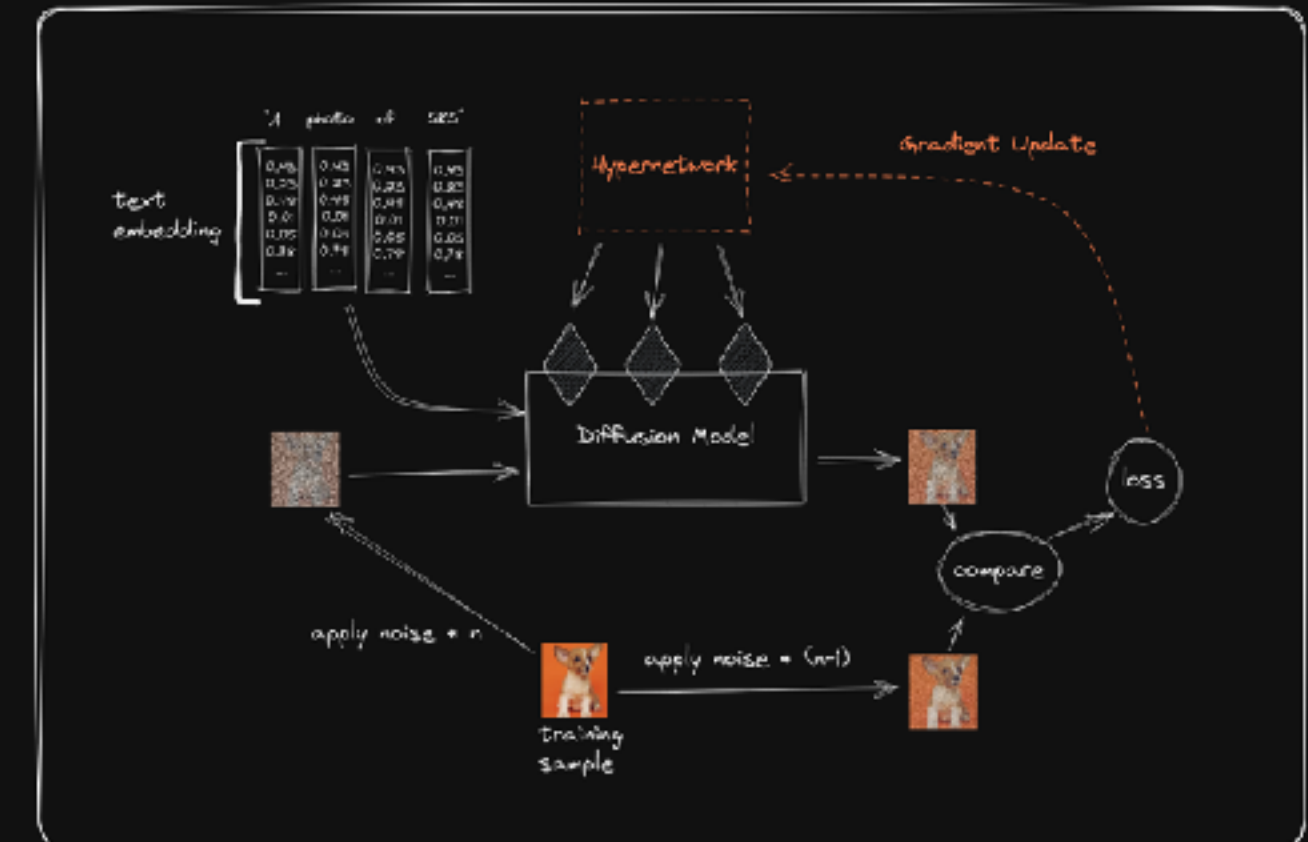
Adds a tiny number of weights to the diffusion model and trains these until the modified model understands the concept

+ quick to train



Hypernetworks

Use a secondary network to predict new weights for the original network. The new weights are swapped in at inference.



Finomhangolás

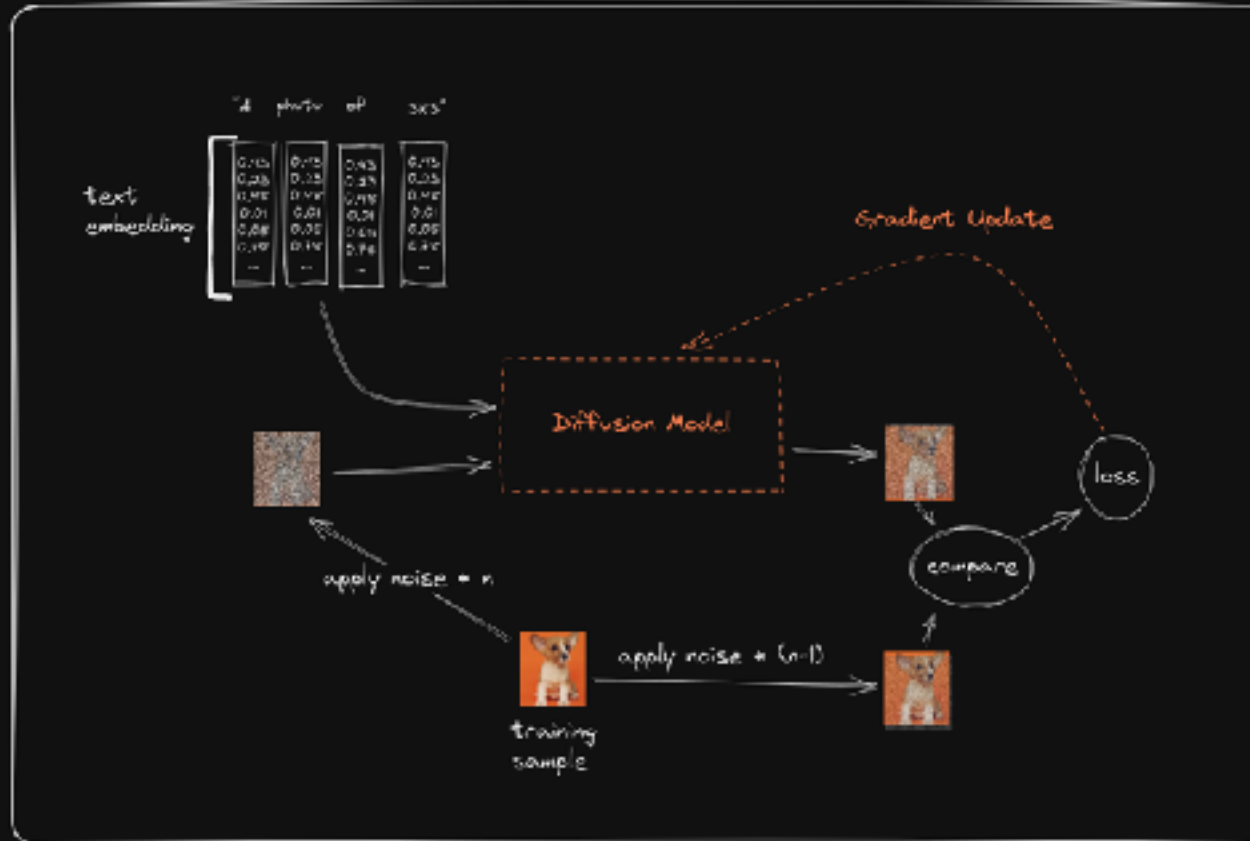
Áttekintés



Dreambooth

Fine-tunes the diffusion model itself until it understands the new concept

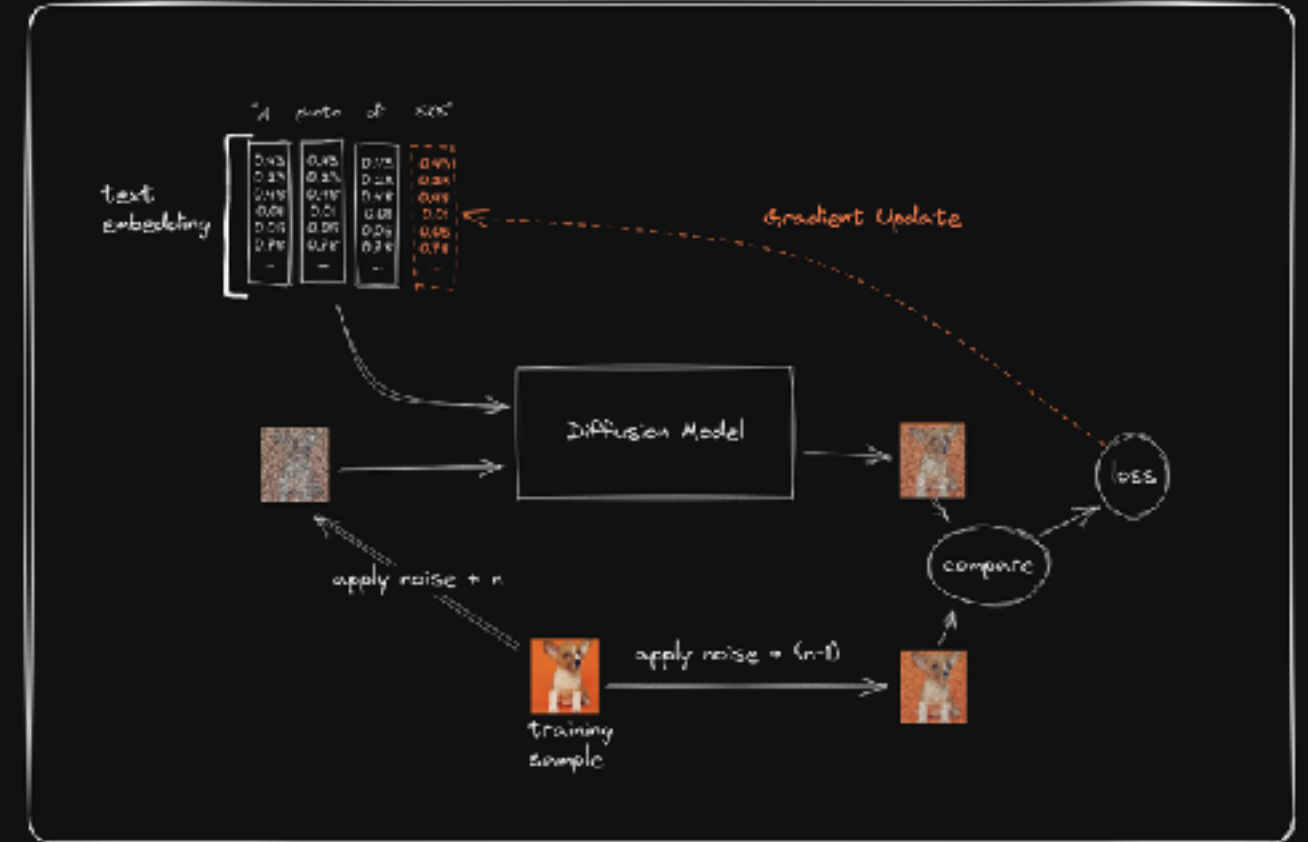
+ probably the most effective
- storage inefficient (whole new model to deal with)



Textual Inversion

Creates a special word embedding which captures the new concept

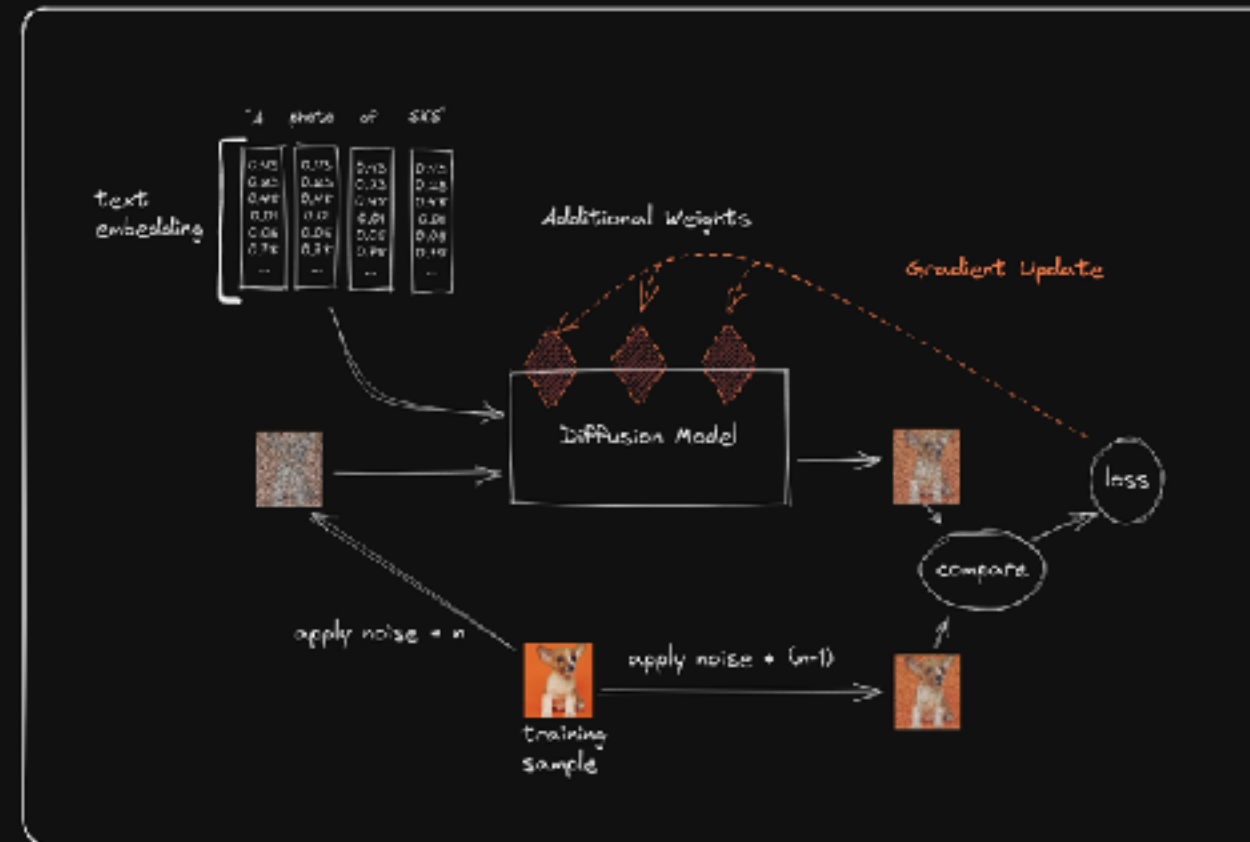
+ output is a tiny embedding



LoRA

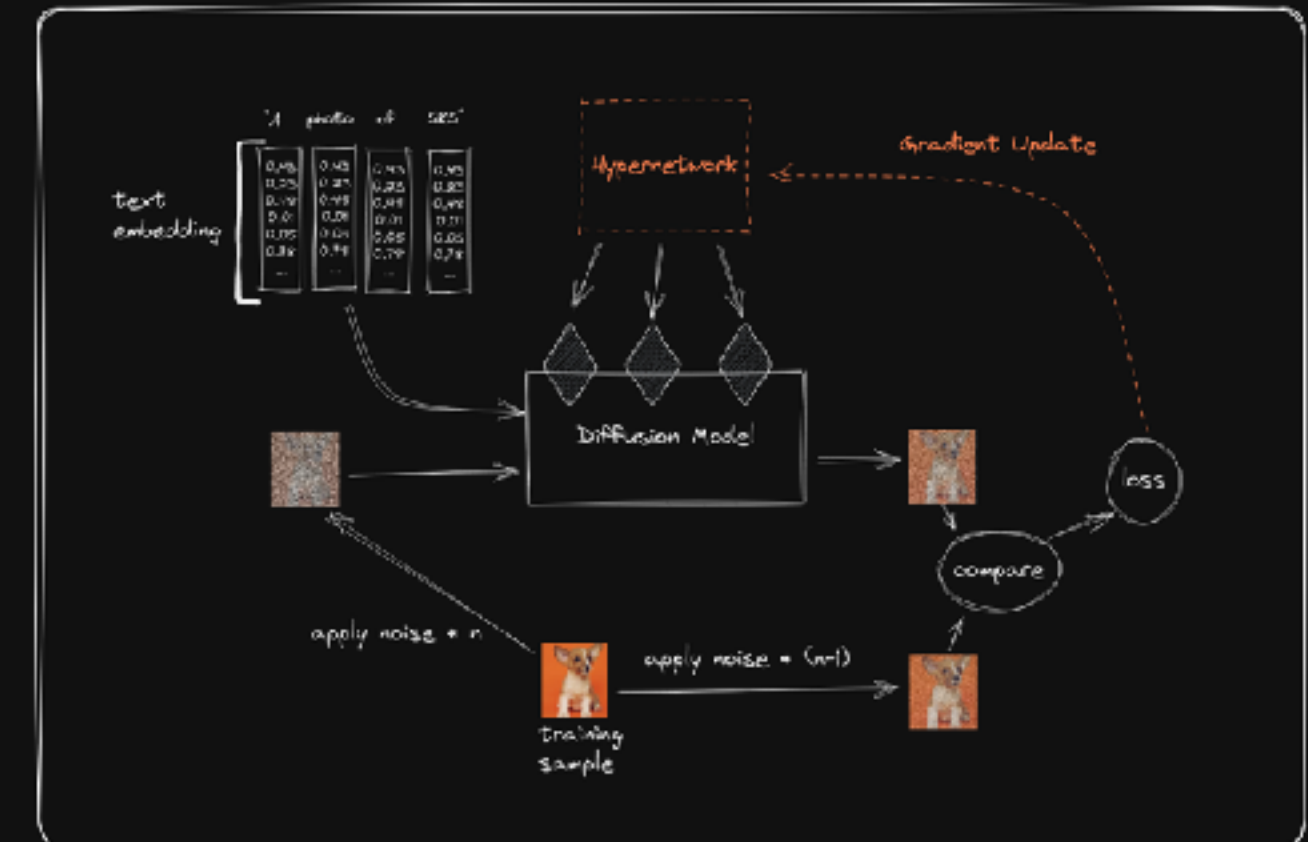
Adds a tiny number of weights to the diffusion model and trains these until the modified model understands the concept

+ quick to train



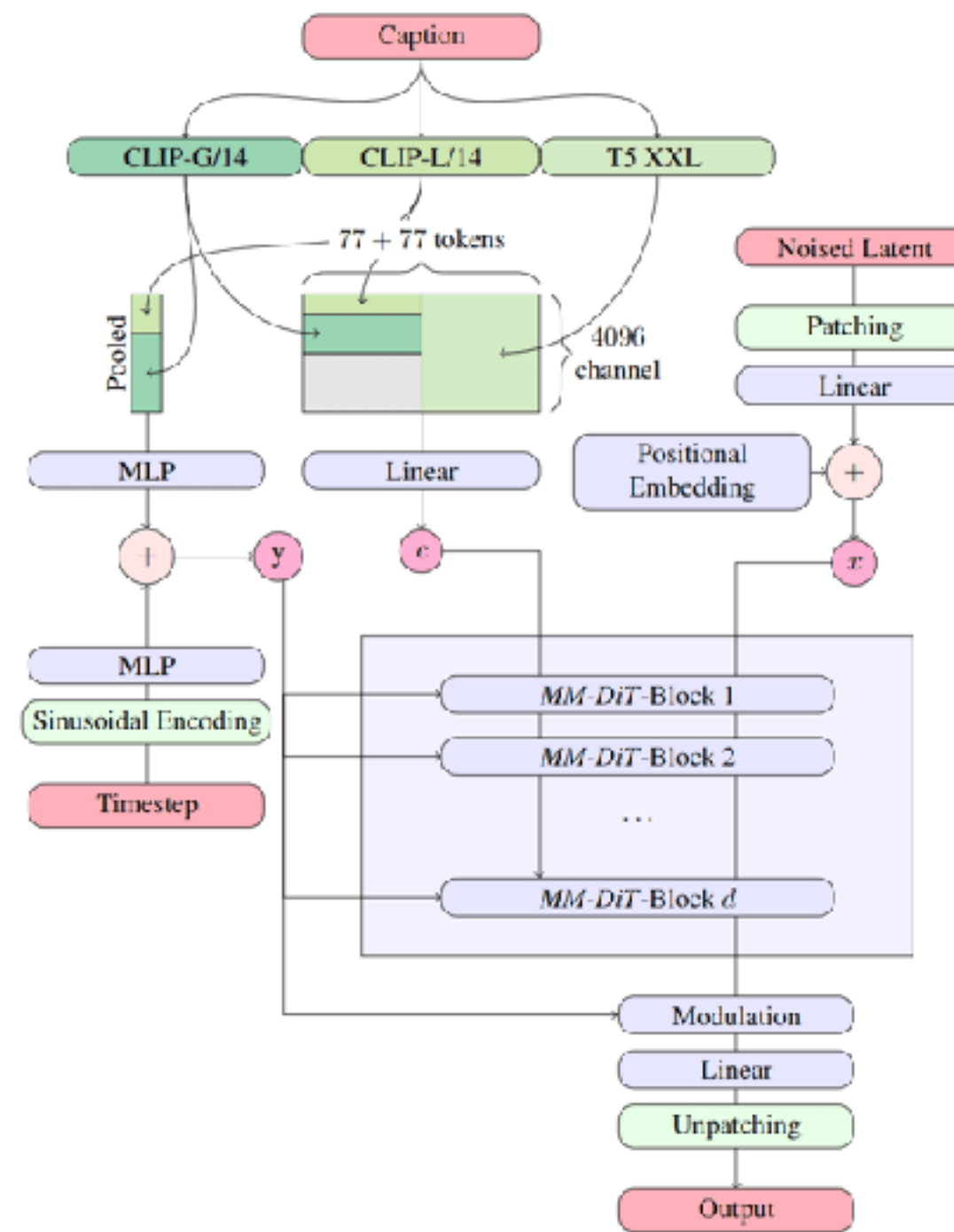
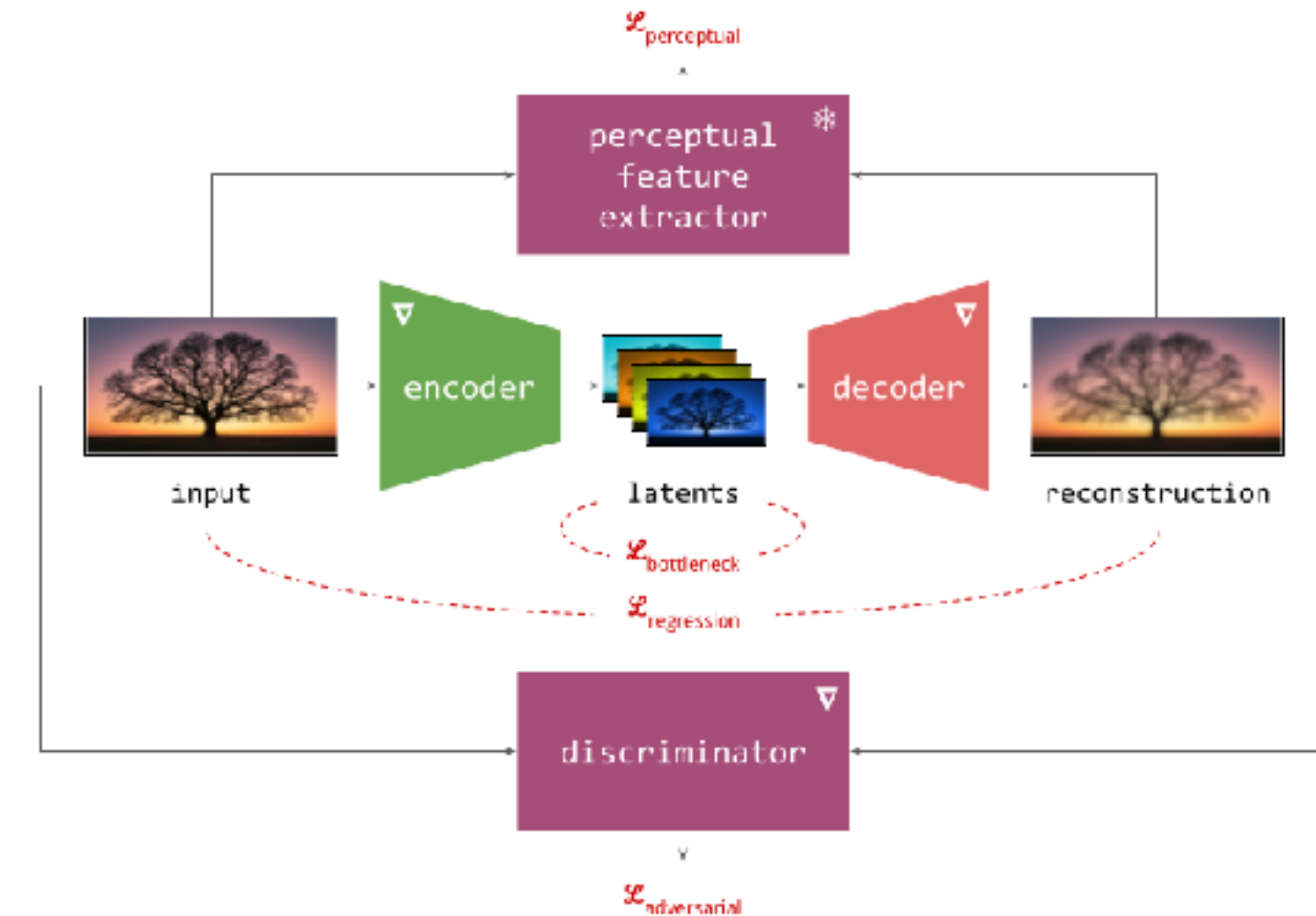
Hypernetworks

Use a secondary network to predict new weights for the original network. The new weights are swapped in at inference.

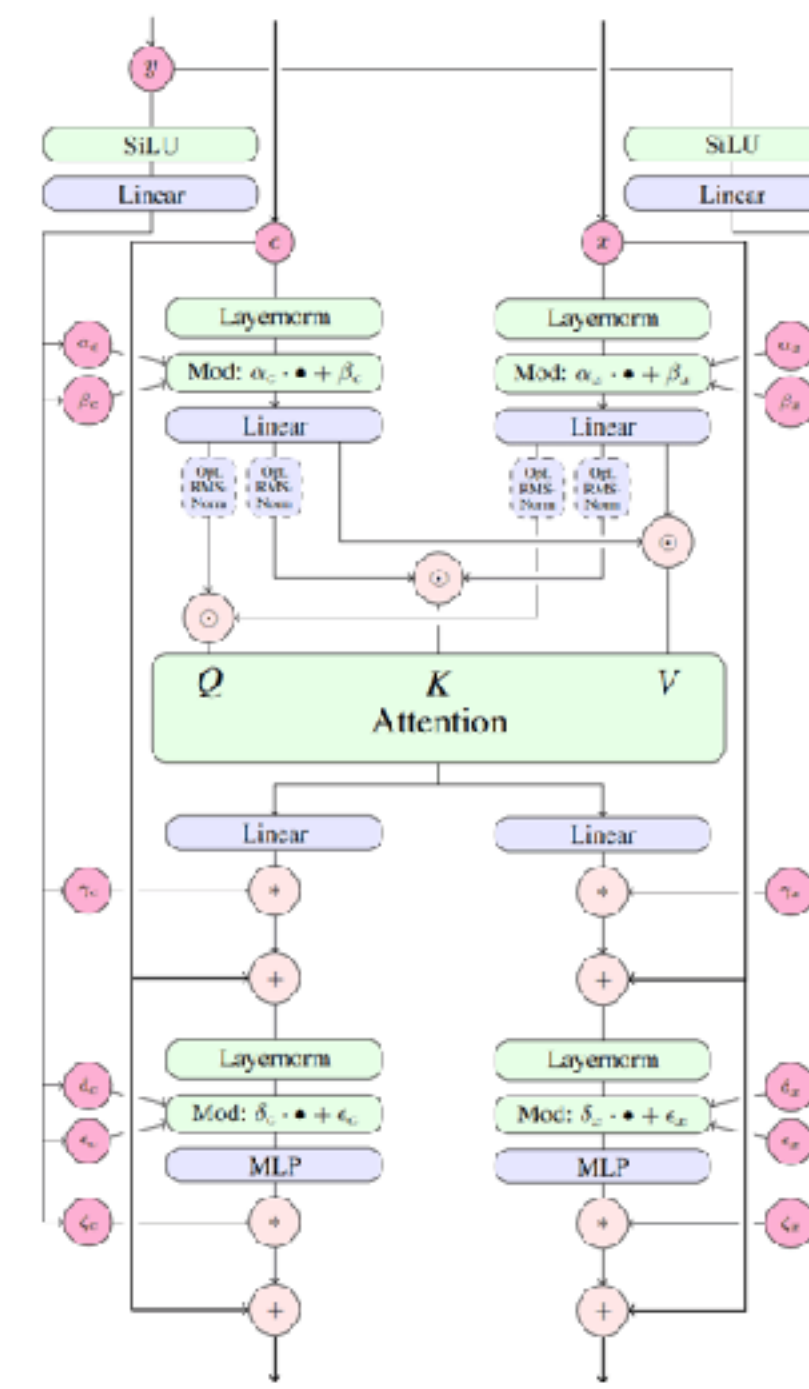


Következő előadás: Diffúziós Generálás a Gyakorlatban

- Látens diffúzió
- Esettanulmányok:
 - Dall-E, Imagen
 - Stable Diffusion 1-3
 - Stb.



(a) Overview of all components.



(b) One MM-DiT block