

# 與影像有關的生成式 AI



【生成式AI導論 2024】第15講：為什麼語言模型用文字接龍，圖片生成不用像素接龍呢？— 淺談生成式人工智慧的生成策略

<https://youtu.be/QbwQR9sjWbs?si=4svPnXaoD1rpcfgv>

# 圖片是由像素所構成

- 影片是由一張一張圖片所構成



video



frame 1



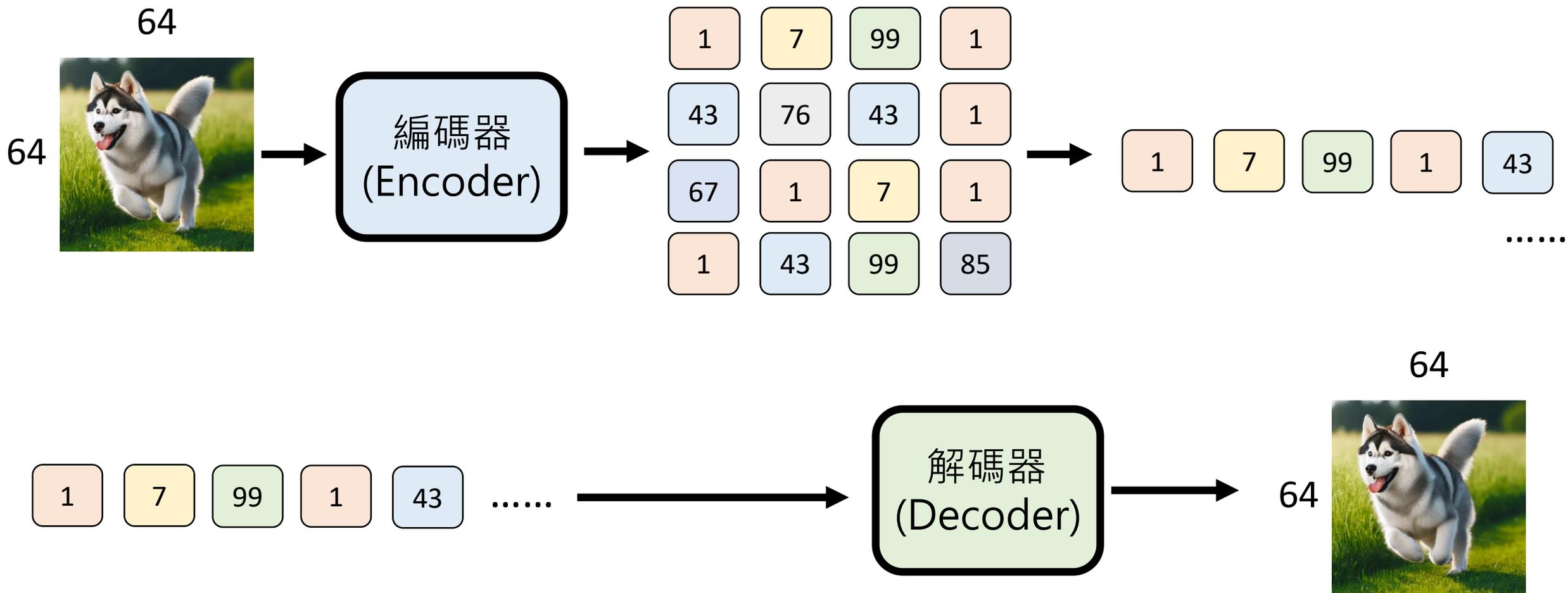
frame 2



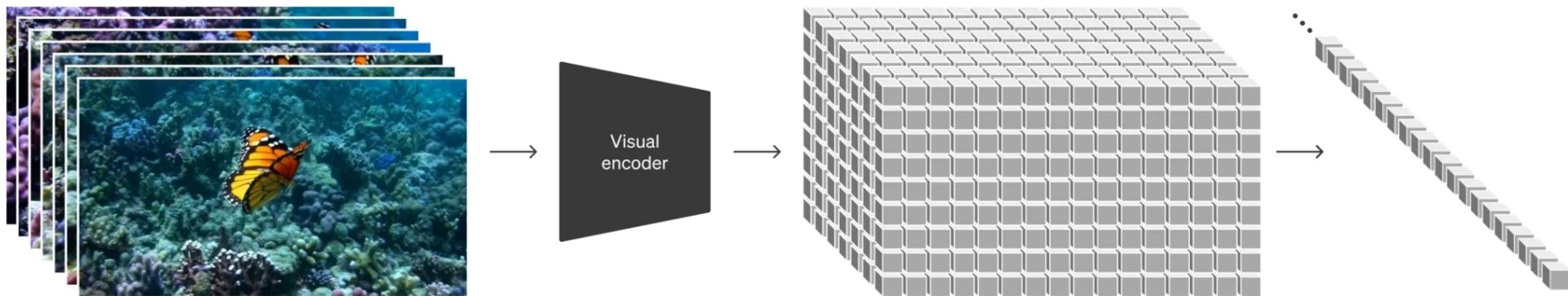
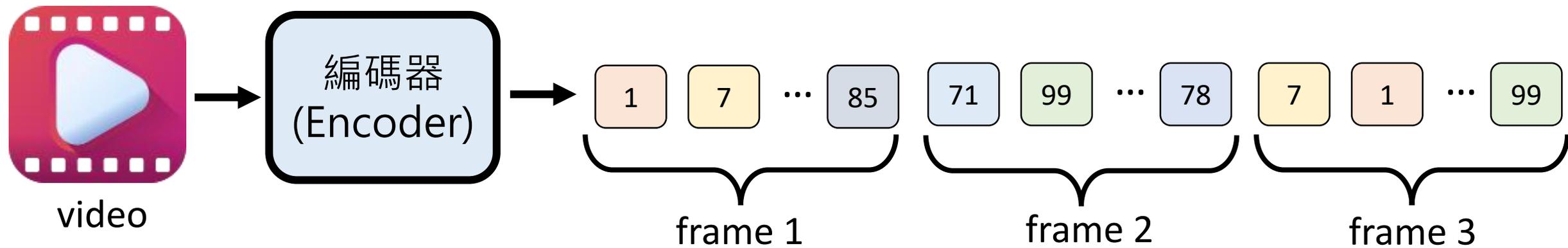
frame 3



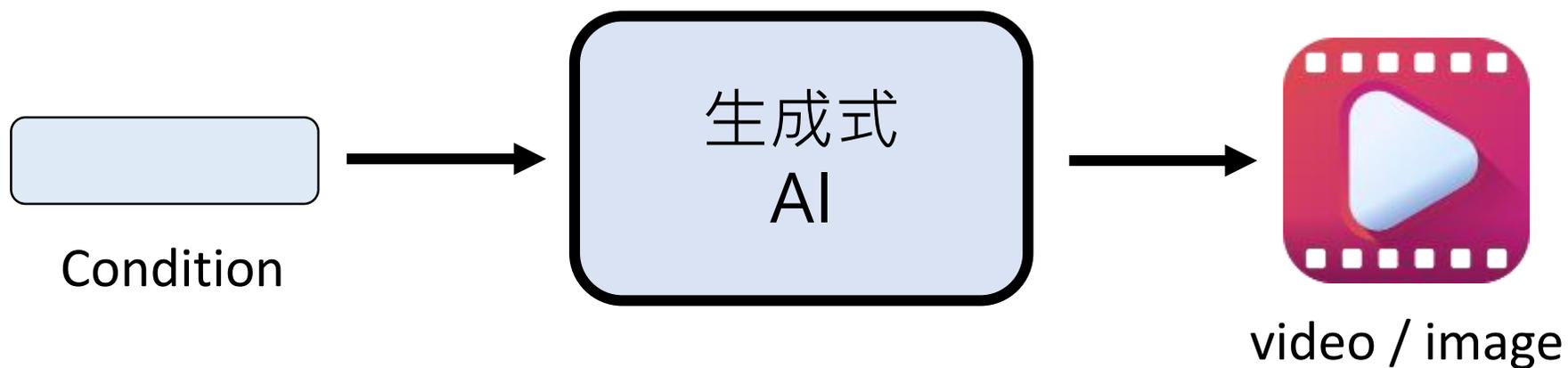
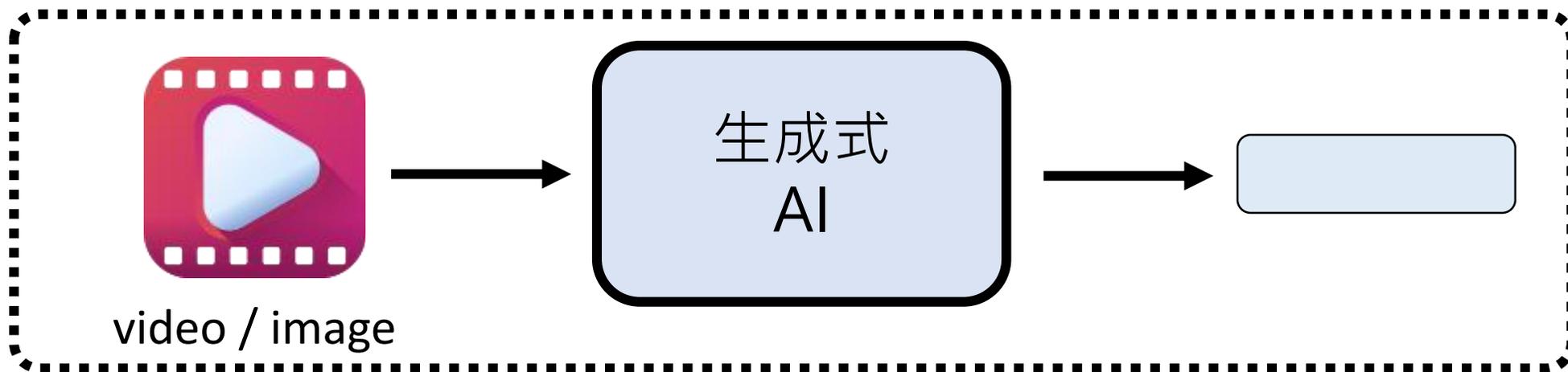
# 今日的人工智慧如何看影像 - 圖片



# 今日的人工智慧如何看影像 – 影片

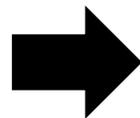


# 與影像有關的生成式AI



GPT-4o

這兩個人人都叫李宏毅



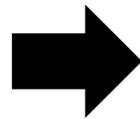
你覺得哪一張圖片中的人比較帥

猜猜看他們分別從事甚麼職業



猜測左邊那位穿著牛仔外套的年輕人可能從事時尚、娛樂或創意領域的工作，例如模特、演員或藝術家。右邊那位穿著白襯衫戴眼鏡的男士看起來更像是專業人士，可能從事學術、技術或商業領域的工作，例如教授、工程師或經理。你能告訴我他們實際從事的職業嗎？

這兩個人都叫李宏毅



你覺得哪一張圖片中的人比較帥

GPT-4o

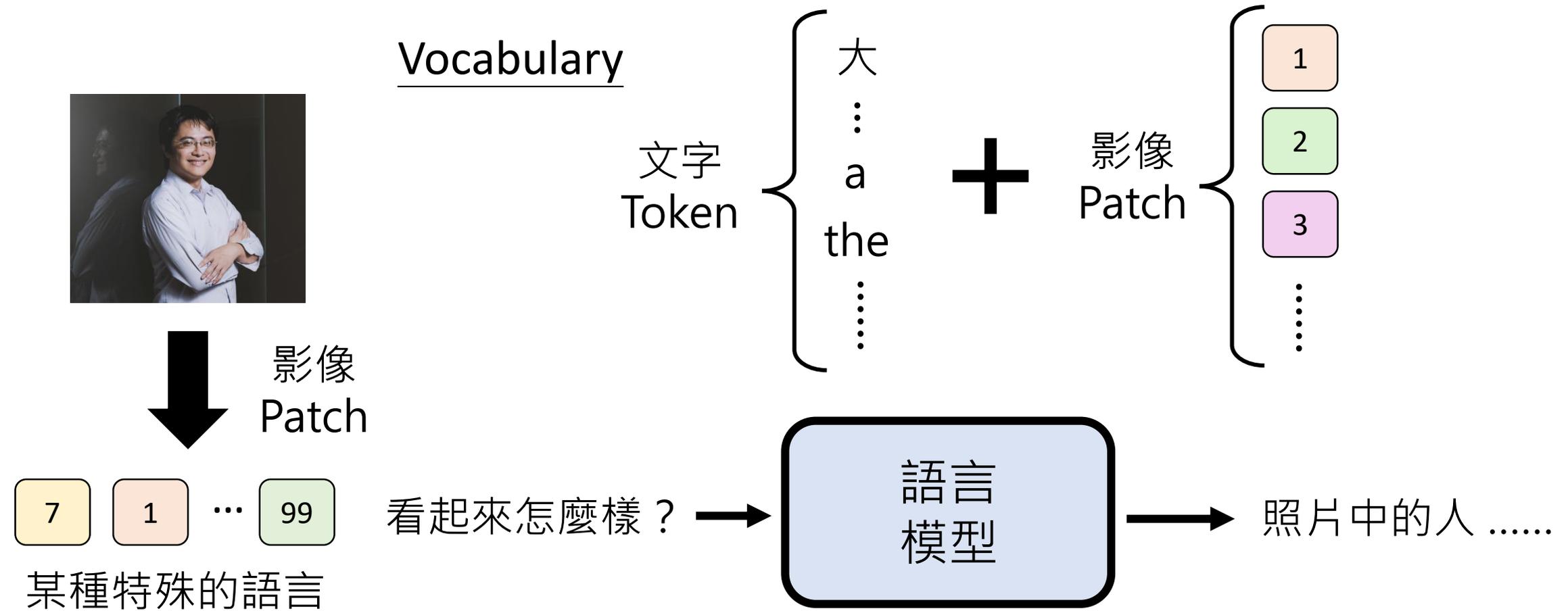


圖片中的兩位風格不同，各有特色。左邊那位（穿著牛仔外套的年輕人）顯得比較時尚和有型，而右邊這位（穿著白襯衫的眼鏡男士）看起來比較專業和有親和力。具體誰比較帥，這很主觀，取決於個人喜好。你覺得呢？

你一定要決定誰比較帥

(讓模型看影像有很多方法，此處只說明了某一種可能性)

# 以影像作為 condition



# 延伸閱讀

## An Introduction to Vision-Language Modeling

Florian Bordes<sup>\*</sup>, Richard Yuanzhe Pang<sup>\*^</sup>, Anurag Ajay<sup>\*\*\*</sup>, Alexander C. Li<sup>\*\*</sup>, Adrien Bardes<sup>\*</sup>, Suzanne Petryk<sup>△</sup>, Oscar Mañas<sup>†</sup>, Zhiqiu Lin<sup>◆</sup>, Anas Mahmoud<sup>†</sup>, Bargav Jayaraman<sup>\*</sup>, Mark Ibrahim<sup>\*</sup>, Melissa Hall<sup>\*</sup>, Yunyang Xiong<sup>\*</sup>, Jonathan Lebensold<sup>♡</sup>, Candace Ross<sup>\*</sup>, Srihari Jayakumar<sup>\*</sup>, Chuan Guo<sup>\*</sup>, Diane Bouchacourt<sup>\*</sup>, Haider Al-Tahan<sup>\*</sup>, Karthik Padthe<sup>\*</sup>, Vasu Sharma<sup>\*</sup>, Hu Xu<sup>\*</sup>, Xiaoqing Ellen Tan<sup>\*</sup>, Megan Richards<sup>\*</sup>, Samuel Lavoie<sup>†</sup>, Pietro Astolfi<sup>\*</sup>, Reyhane Askari Hemmat<sup>\*</sup>, Jun Chen<sup>\*\*\*◇</sup>, Kushal Tirumala<sup>\*</sup>, Rim Assouel<sup>†</sup>, Mazda Moayeri<sup>▽</sup>, Arjang Talattof<sup>\*</sup>, Kamalika Chaudhuri<sup>\*</sup>, Zechun Liu<sup>\*</sup>, Xilun Chen<sup>\*</sup>, Quentin Garrido<sup>\*</sup>, Karen Ullrich<sup>\*</sup>, Aishwarya Agrawal<sup>†\*</sup>, Kate Saenko<sup>\*</sup>, Asli Celikyilmaz<sup>\*</sup> and Vikas Chandra<sup>\*</sup>

<sup>\*</sup>Meta

<sup>\*\*</sup>Work done while at Meta

<sup>†</sup>Université de Montréal, Mila

<sup>♡</sup>McGill University, Mila

<sup>†</sup>University of Toronto

<sup>◆</sup>Carnegie Mellon University

<sup>◆</sup>Massachusetts Institute of Technology

<sup>^</sup>New York University

<sup>△</sup>University of California, Berkeley

<sup>▽</sup>University of Maryland

<sup>◇</sup>King Abdullah University of Science and Technology

<sup>\*</sup>Canada CIFAR AI Chair

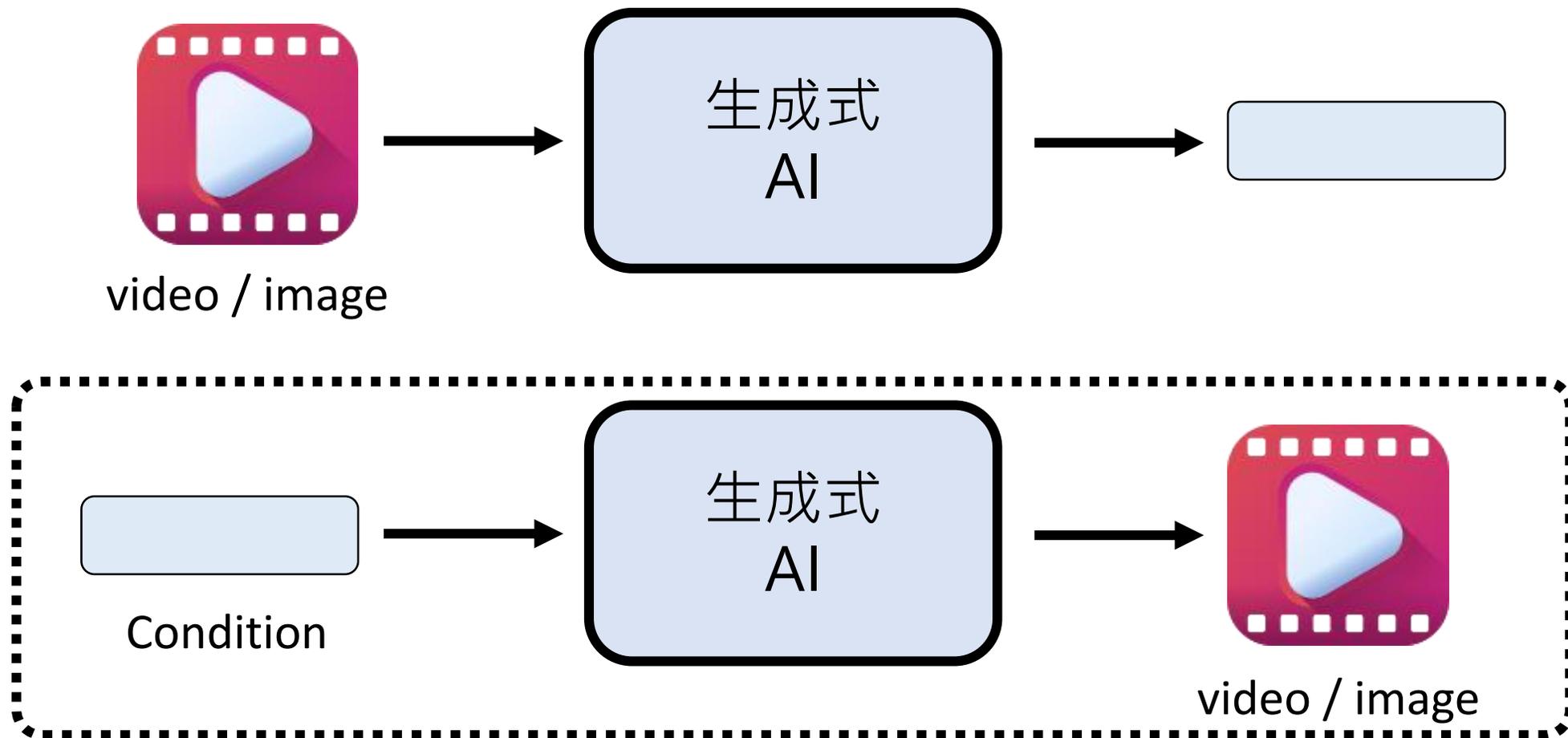
Core contributors, random ordering

Additional contributors, random ordering

Senior contributors, random ordering

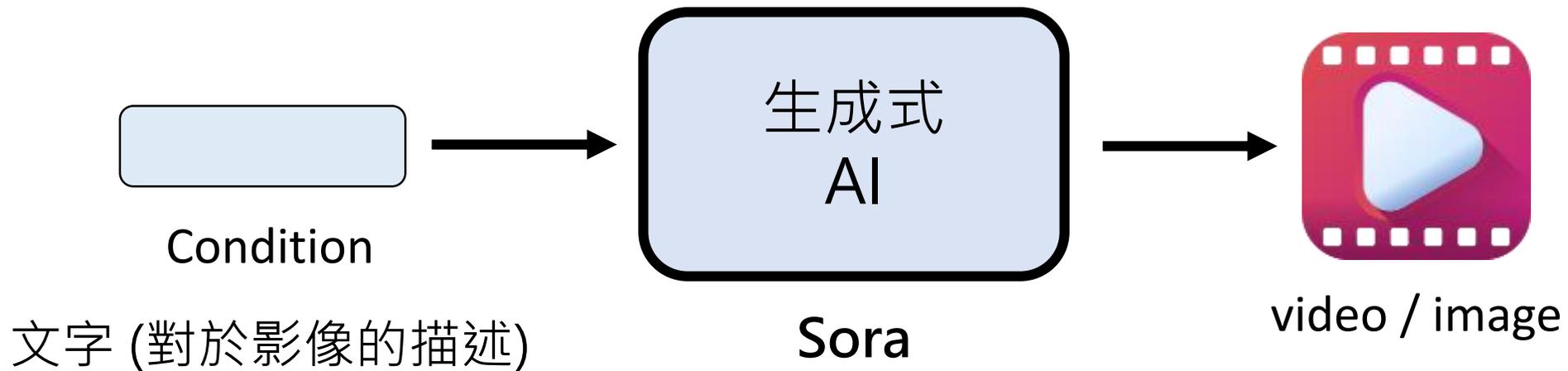
<https://arxiv.org/abs/2405.17247>

# 與影像有關的生成式AI



# 生成影像的生成式 AI — 文字生影像

- 以文字作為 condition



Sora 並沒有真的開放使用

<https://openai.com/sora>

<https://openai.com/research/video-generation-models-as-world-simulators>

# 生成影像的生成式 AI — 文字生影像 Sora



Animated scene features a close-up of a short fluffy monster kneeling beside a melting red candle. The art style is 3D and realistic, with a focus on lighting and texture. ....



New York City submerged like Atlantis. Fish, whales, sea turtles and sharks swim through the streets of New York.

# 生成影像的生成式 AI — 文字生影像 Sora

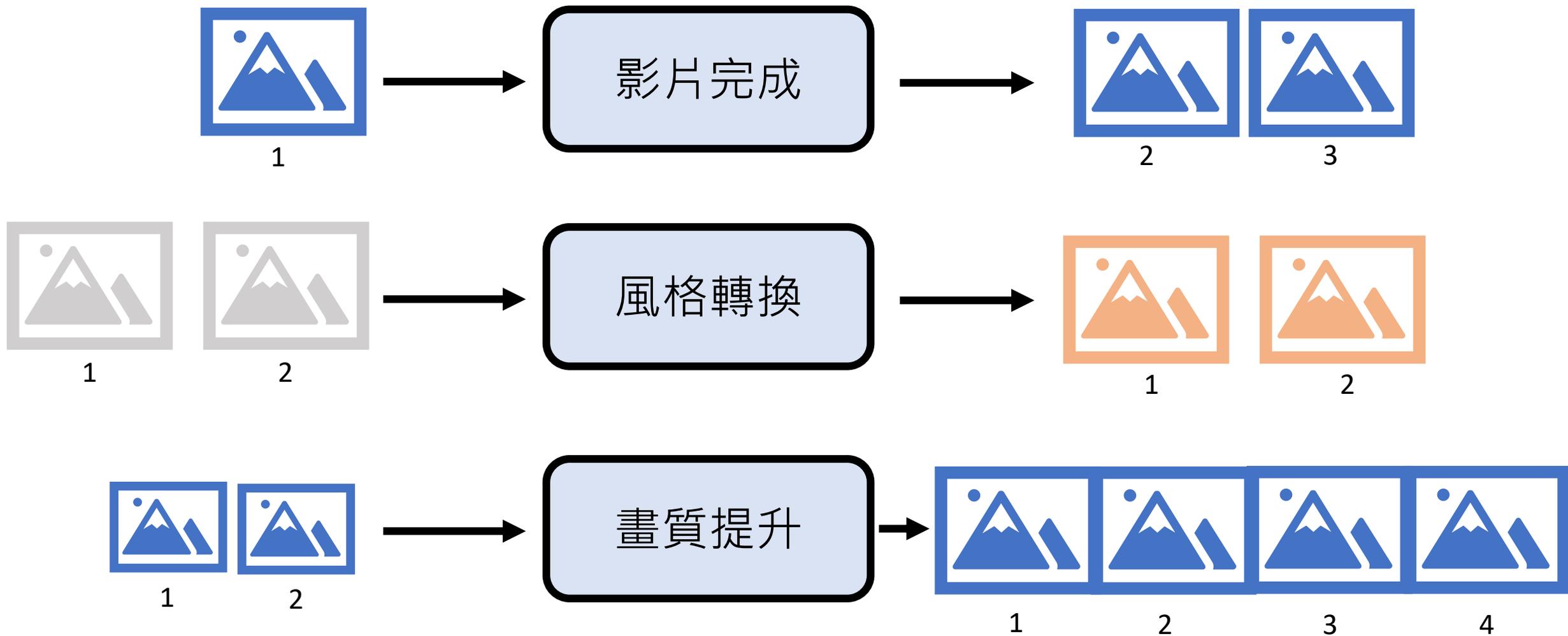


Five gray wolf pups frolicking and chasing each other around a remote gravel road, surrounded by grass. The pups run and leap, chasing each other, and nipping at each other, playing.



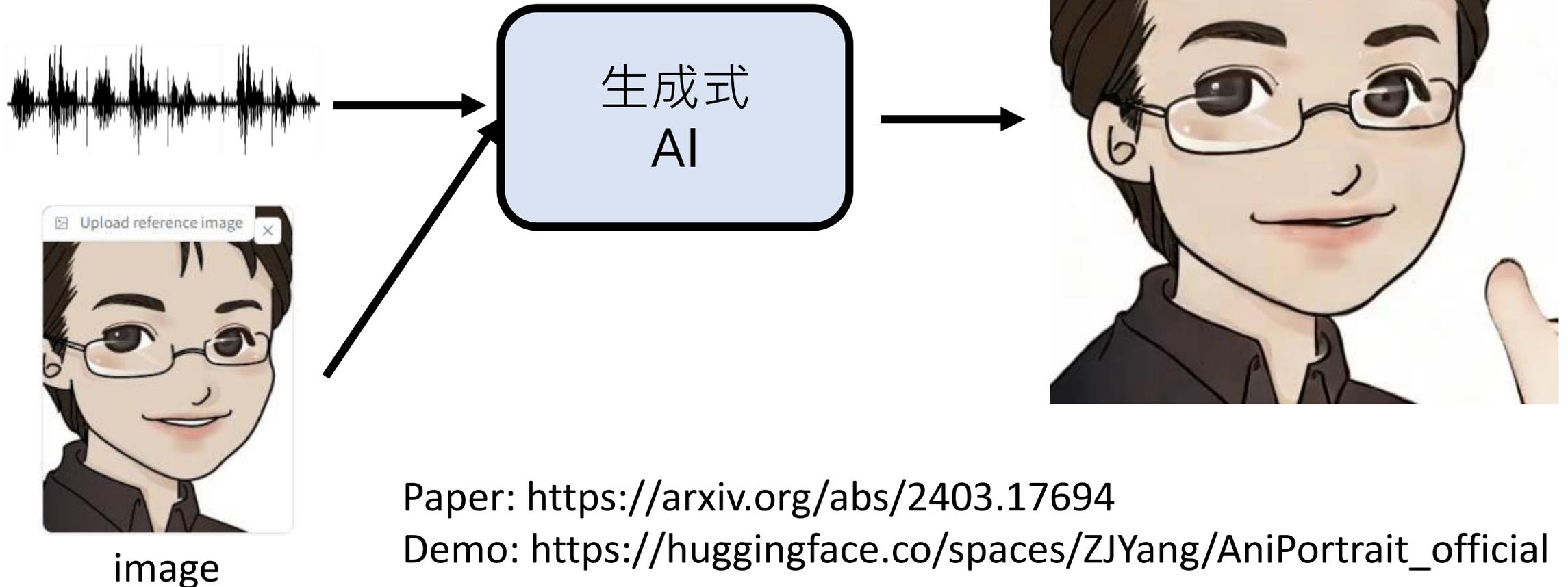
Archeologists discover a generic plastic chair in the desert, excavating and dusting it with great care.

# 生成影像的生成式 AI — 影像生影像



# 生成影像的生成式 AI — 其他輸入生影像

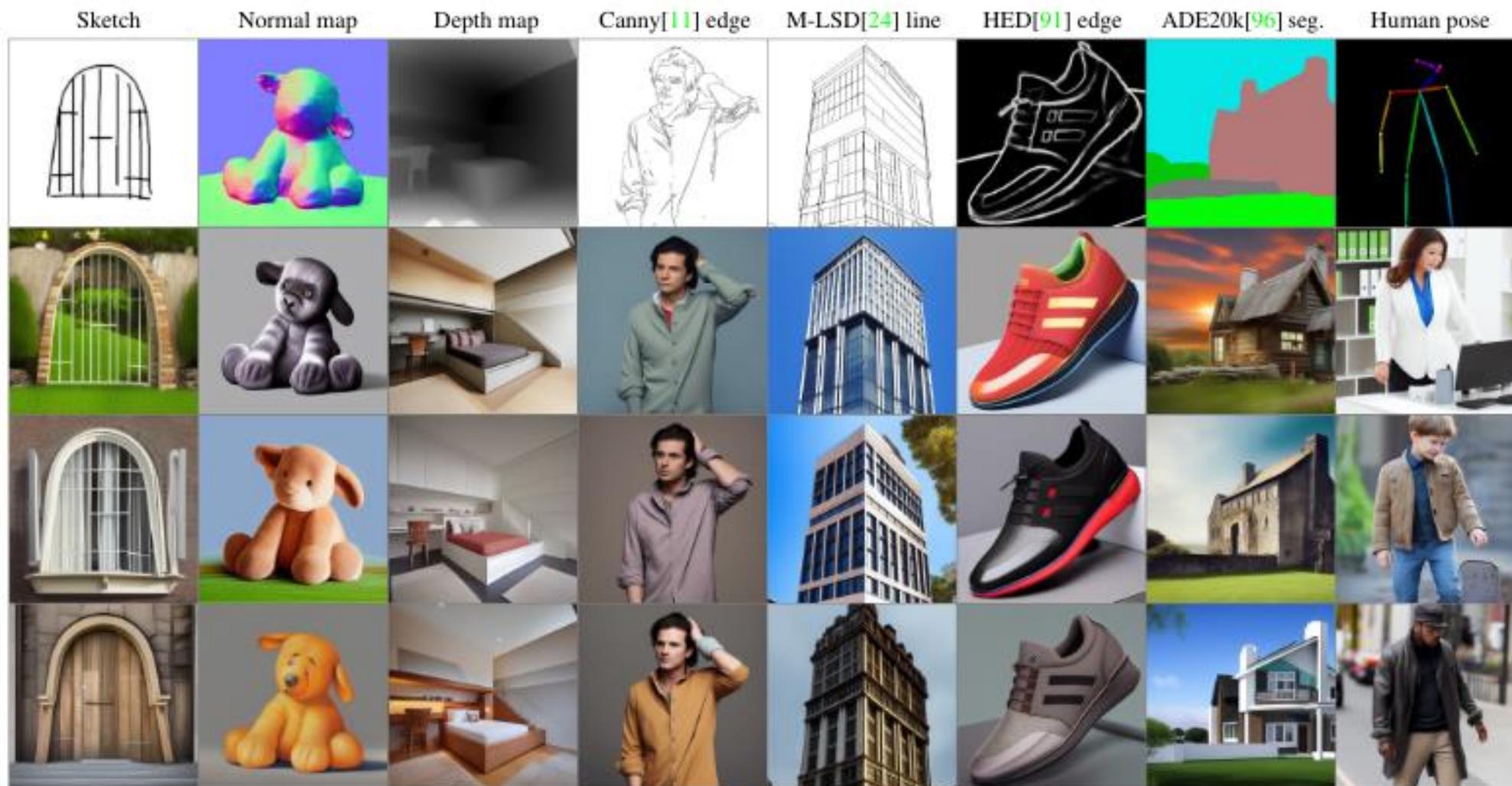
- Talking Head



Paper: <https://arxiv.org/abs/2403.17694>

Demo: [https://huggingface.co/spaces/ZJYang/AniPortrait\\_official](https://huggingface.co/spaces/ZJYang/AniPortrait_official)

# 生成影像的生成式 AI — 其他輸入生影像



# 以文字生圖為例



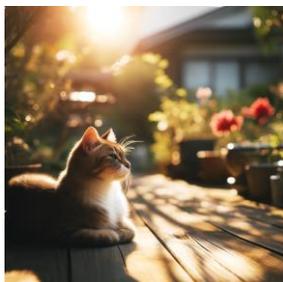
## 訓練資料



一隻在奔跑的狗



雪地裡的貓



陽光下的貓



沙灘上的狗

# 以文字生圖為例

<https://laion.ai/blog/laion-5b/>

**LAION**  
**5.85B**

Backend url:

Index:

Search:  🔍 📷 ⬇️

[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions   
Display full captions   
Display similarities   
Safe mode   
Hide duplicate urls   
Hide (near) duplicate images   
Search over   
Search with multilingual clip



french cat



french cat



How to tell if your feline is french. He wears a b...



イケメン猫モデル「トキ・ナンタケット」がかっこいい - NAVER まとめ



Hilarious pics of funny cats! funnycatsgif.com



網友挑戰「加幾筆畫出最創意貓咪圖片」, 這隻貓贏了



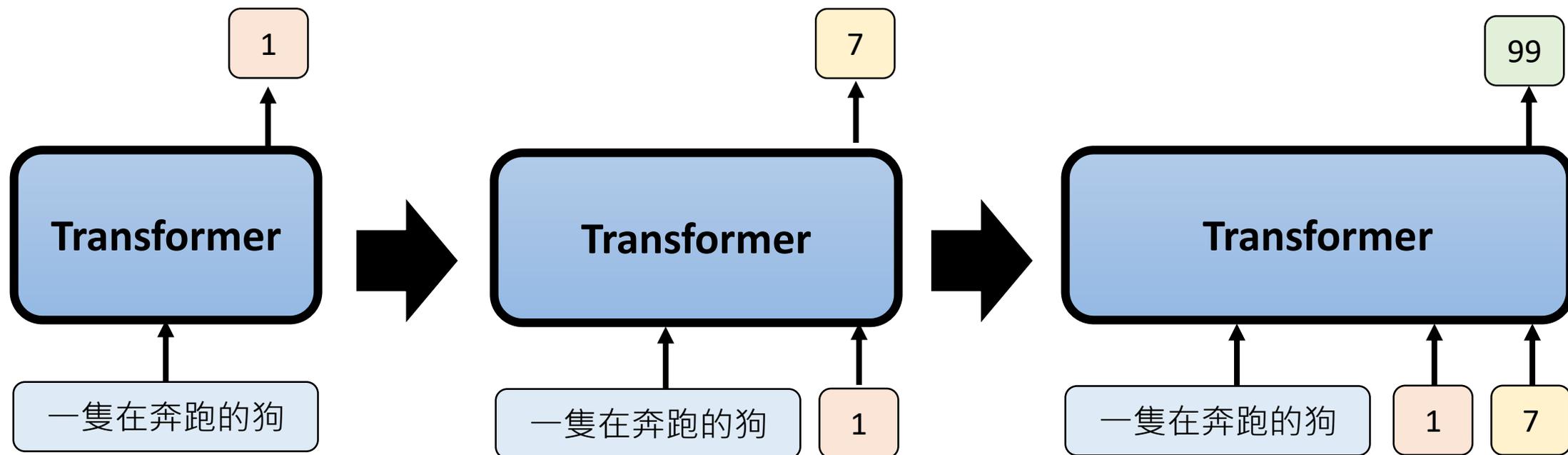
cat in a suit Georgian sells tomatoes



French Bread Cat Loaf

# 以文字生圖為例

一隻在奔跑的狗 1 7 99 1 43 .....

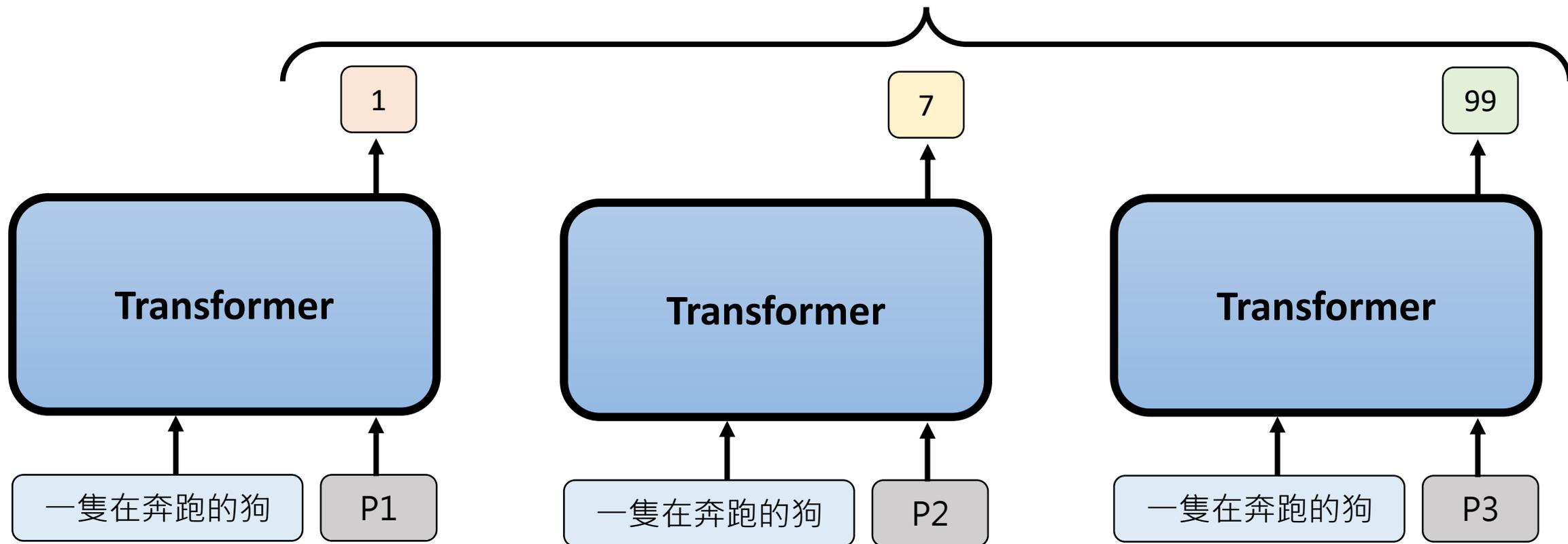


(為了方便同學們理解，本頁投影片對於模型做了大量的簡化)

# 以文字生圖為例

一隻在奔跑的狗 1 7 99 1 43 .....

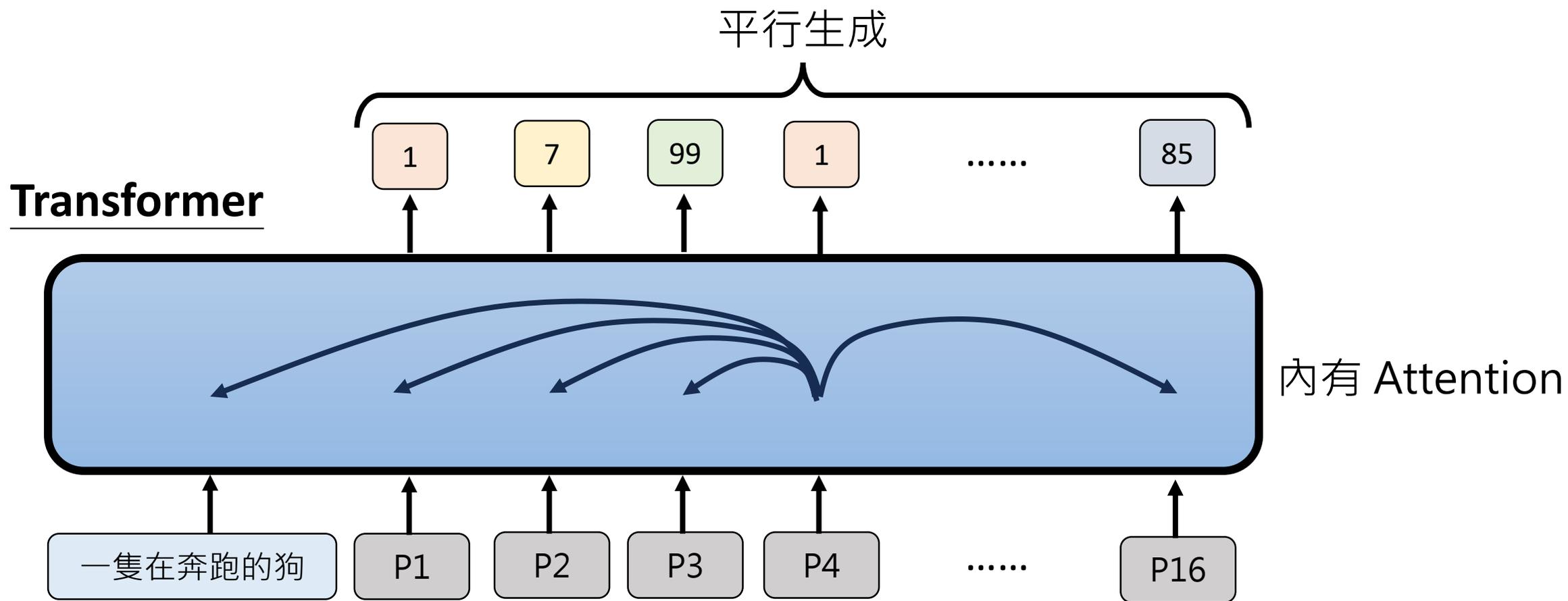
平行生成



(為了方便同學們理解，本頁投影片對於模型做了大量的簡化)

# 以文字生圖為例

最終還是各自獨立生成



(為了方便同學們理解，本頁投影片對於模型做了大量的簡化)

# 如何評量影像生成的好壞

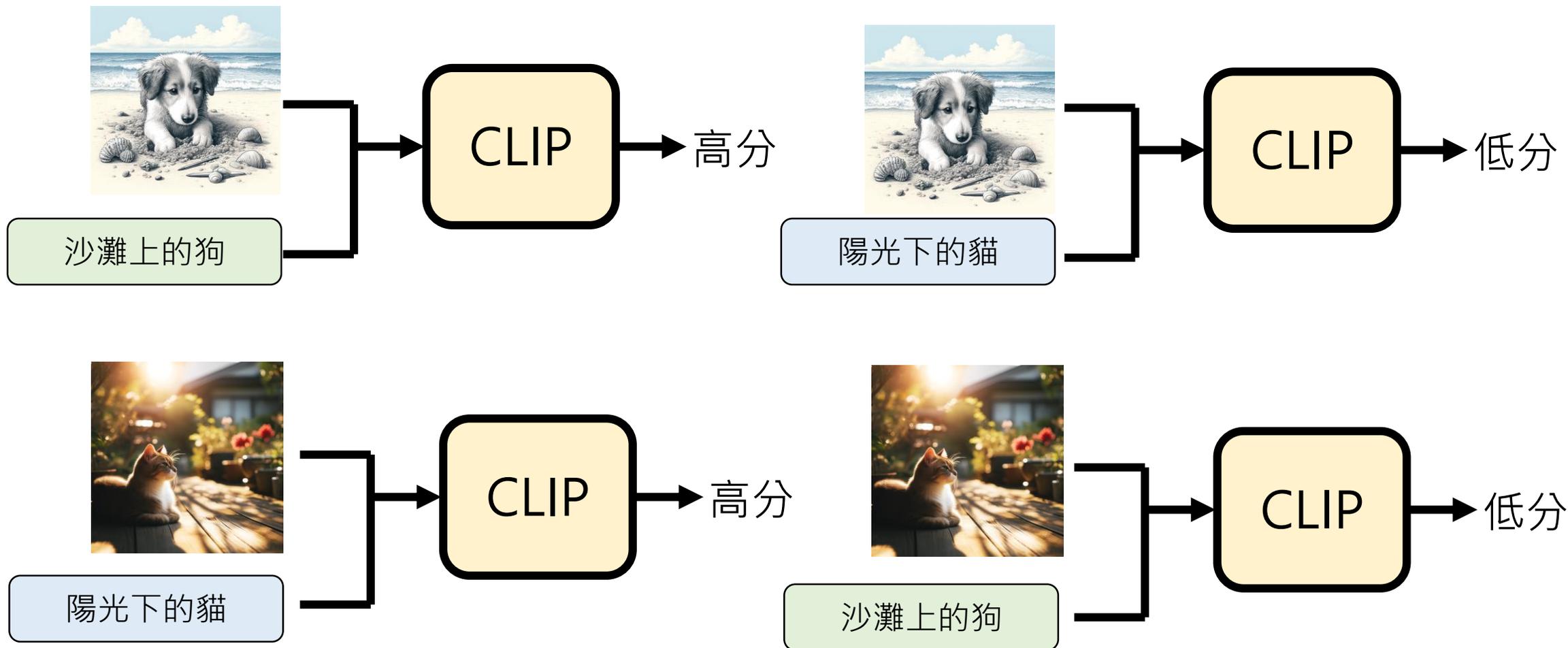
<https://arxiv.org/abs/2103.00020>



怎麼知道生成結果好不好？

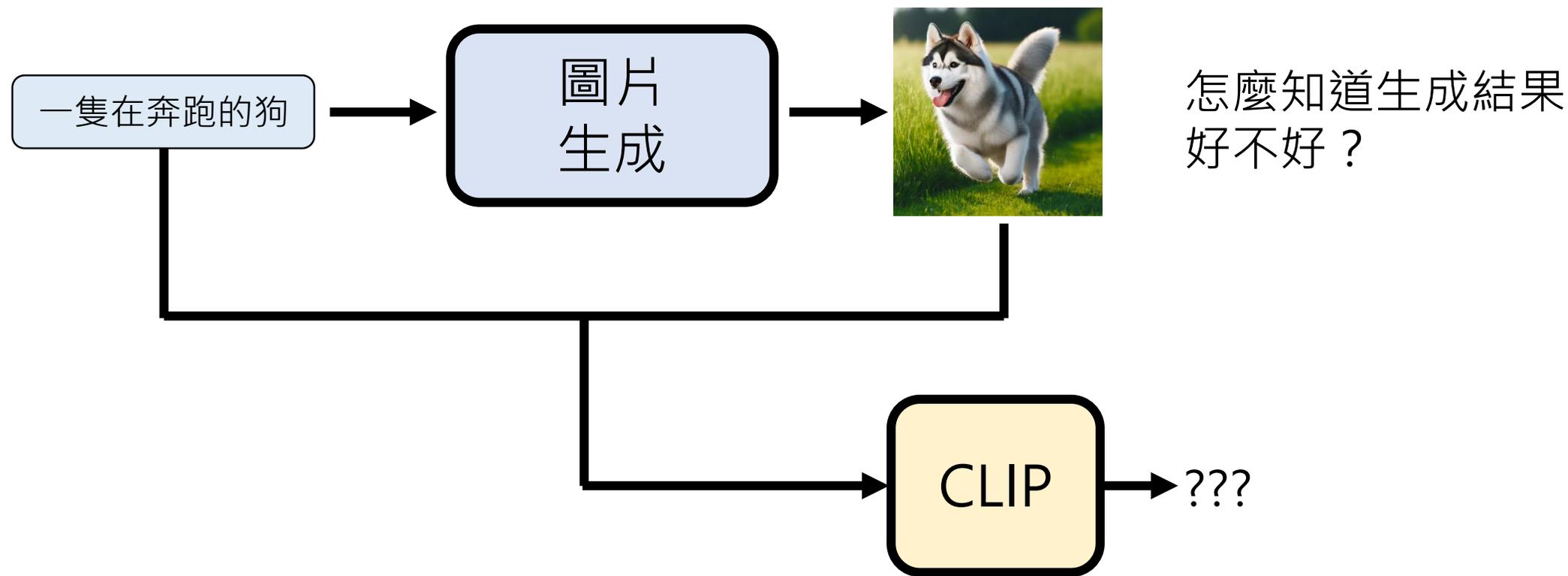
# 如何評量影像生成的好壞

<https://arxiv.org/abs/2103.00020>



# 如何評量影像生成的好壞

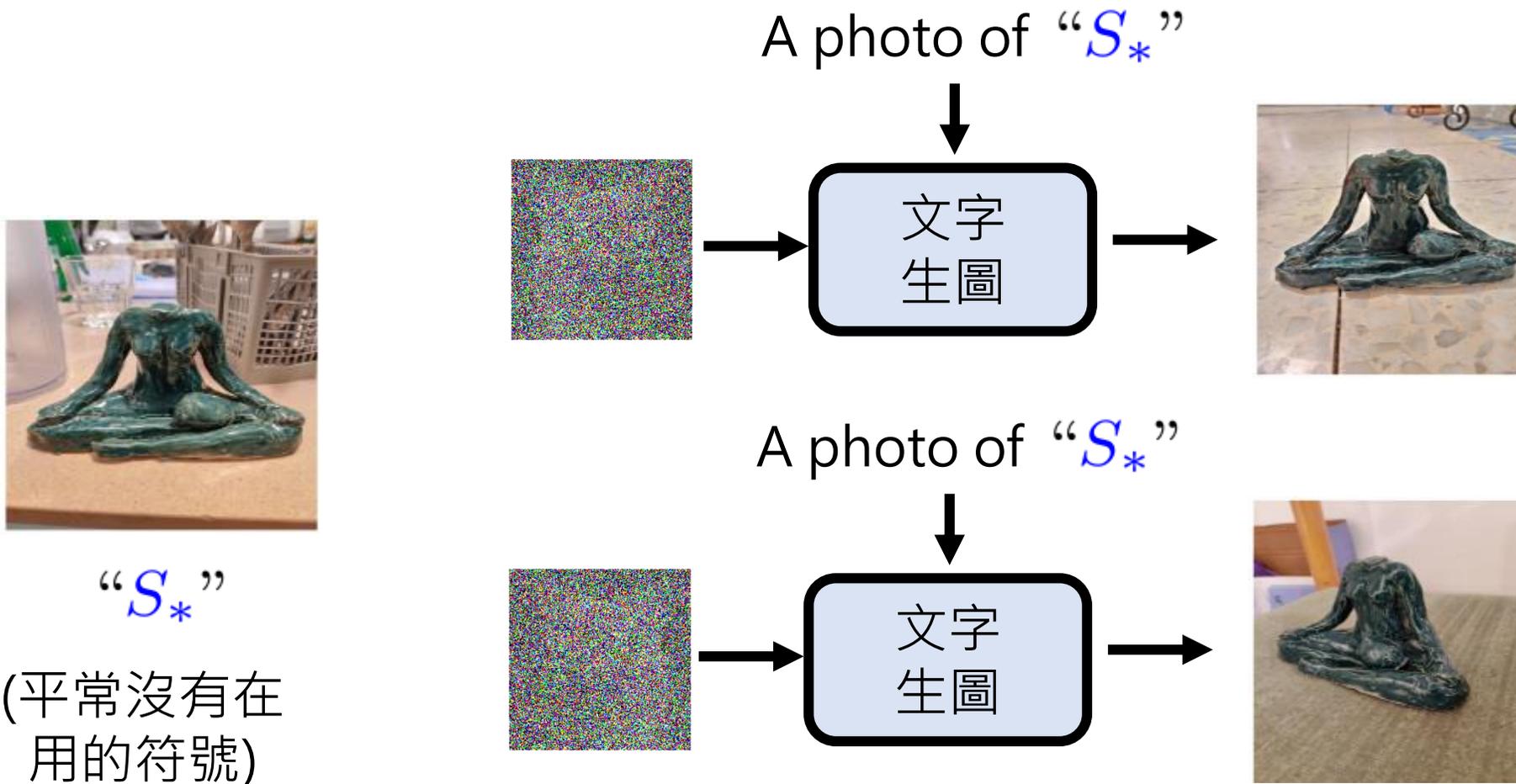
<https://arxiv.org/abs/2103.00020>





# 個人化的圖像生成

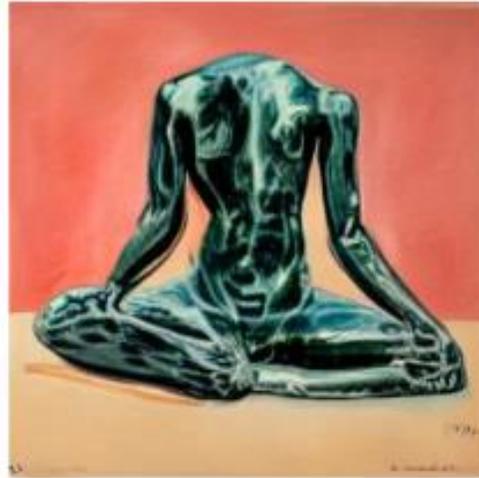
<https://arxiv.org/abs/2208.01618>  
<https://arxiv.org/abs/2208.12242>



Source of image: <https://arxiv.org/abs/2208.01618>



Input samples  $\xrightarrow{\text{invert}}$  “ $S_*$ ”



“An oil painting of  $S_*$ ”



“App icon of  $S_*$ ”

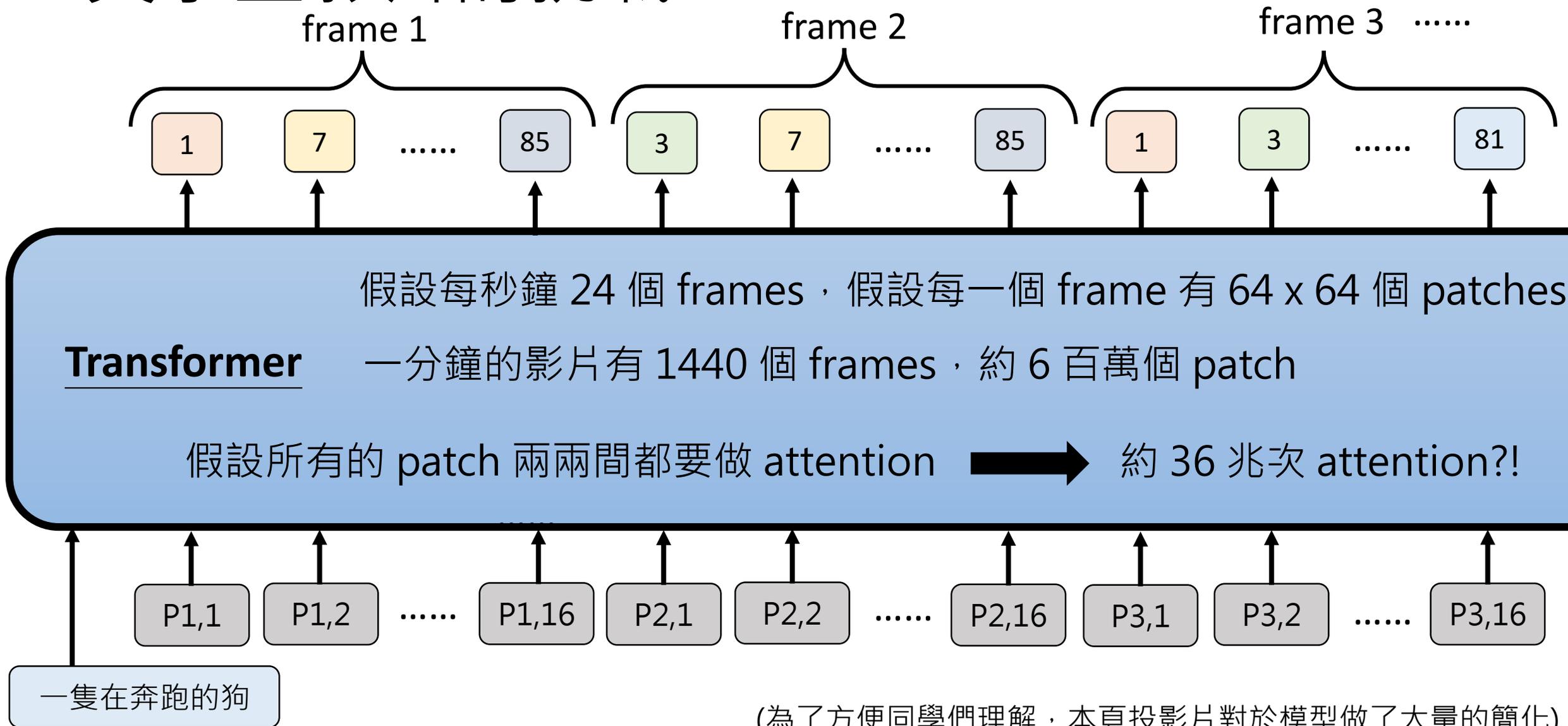


“Elmo sitting in the same pose as  $S_*$ ”

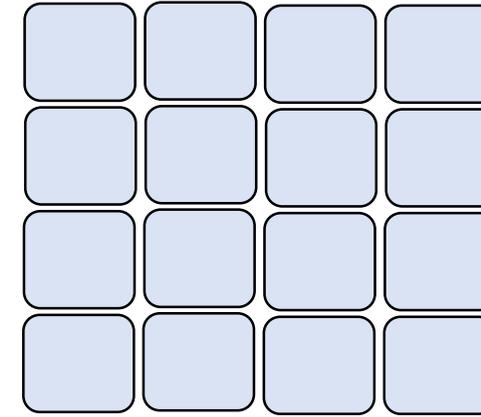
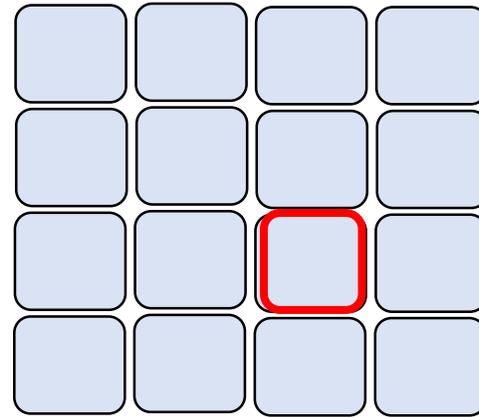
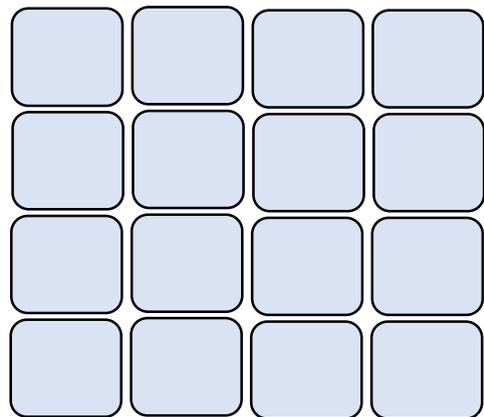


“Crochet  $S_*$ ”

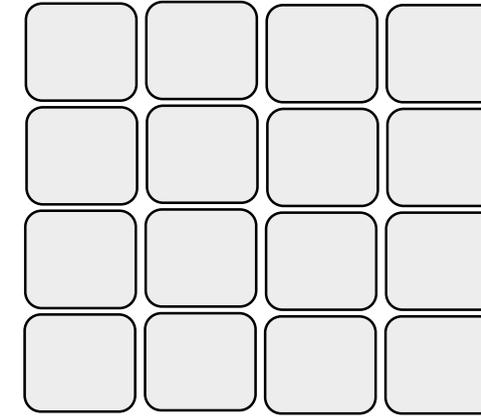
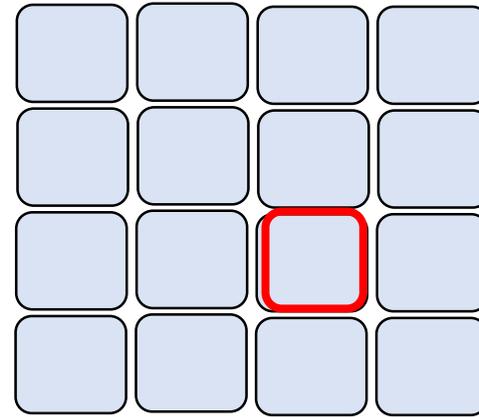
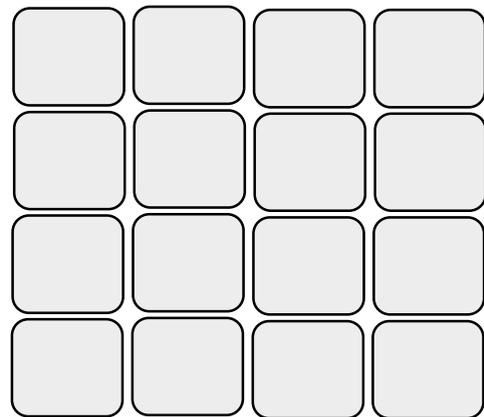
# 文字生影片的挑戰



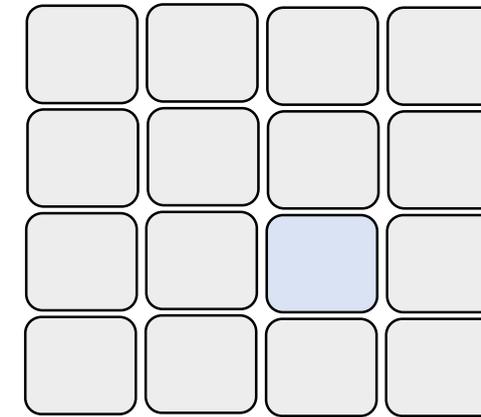
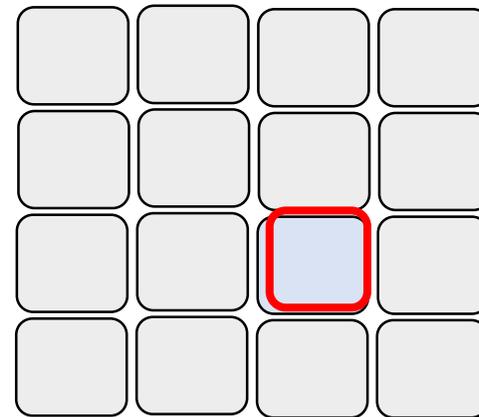
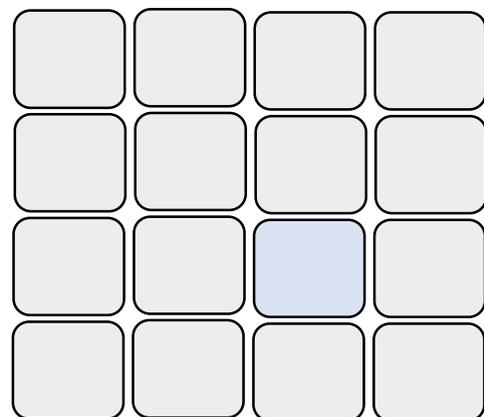
Spatio-temporal  
Attention (3D)

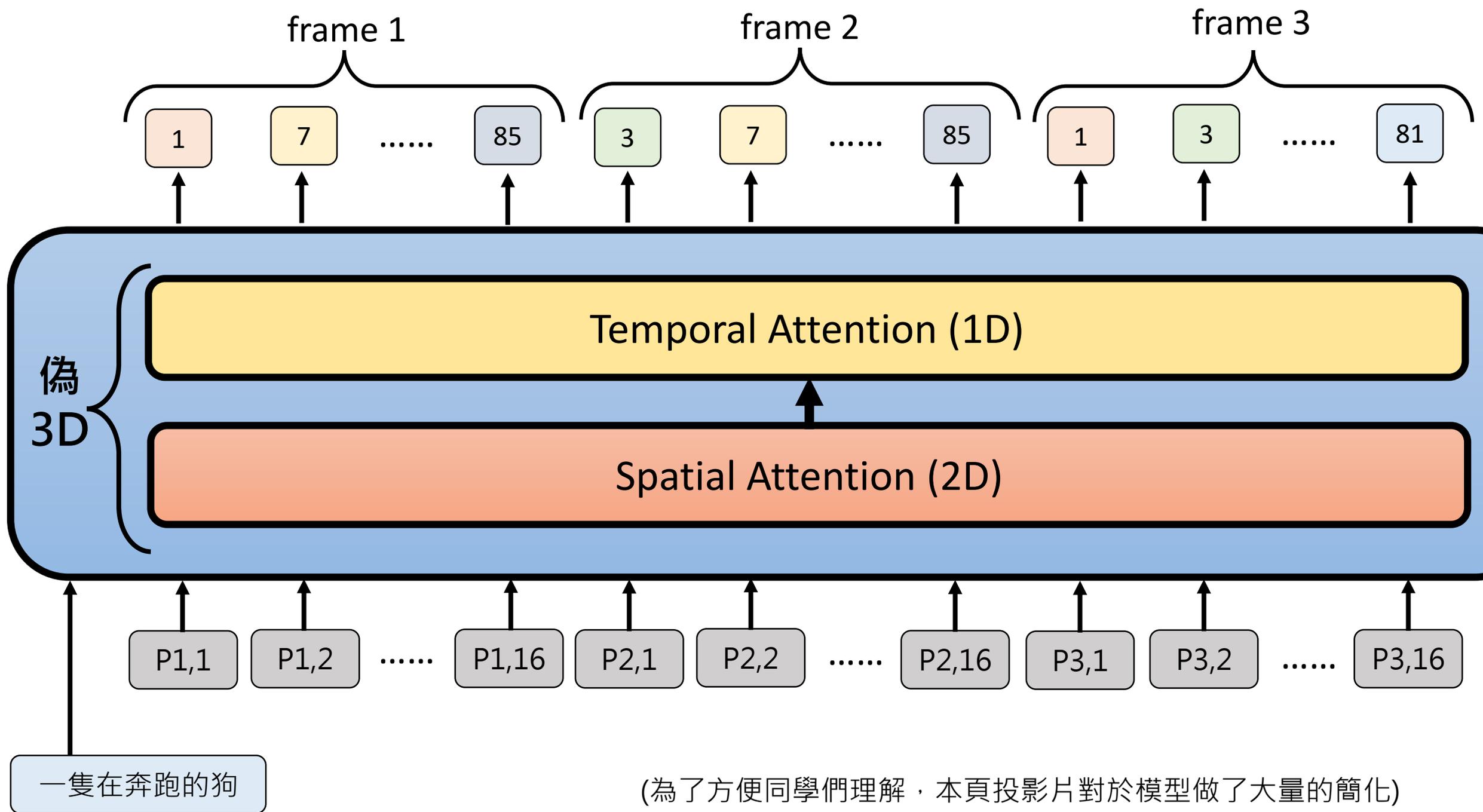


Spatial Attention  
(2D)



Temporal  
Attention (1D)





(為了方便同學們理解，本頁投影片對於模型做了大量的簡化)

假設每秒鐘 24 個 frames，假設每一個 frame 有 64 x 64 個 patches

一分鐘的影片有 1440 個 frames，約 6 百萬個 patch

Spatio-temporal  
Attention (3D)

假設所有的 patch 兩兩間都要做 attention



約 36 兆次 attention?!

Spatial Attention  
(2D)

同一個 frame 中的 patch 才做 attention

$$(64 \times 64)^2 \times 1440$$



約 240 億次 attention

Temporal  
Attention (1D)

不同 frame 中同樣位置的 patch 才做 attention

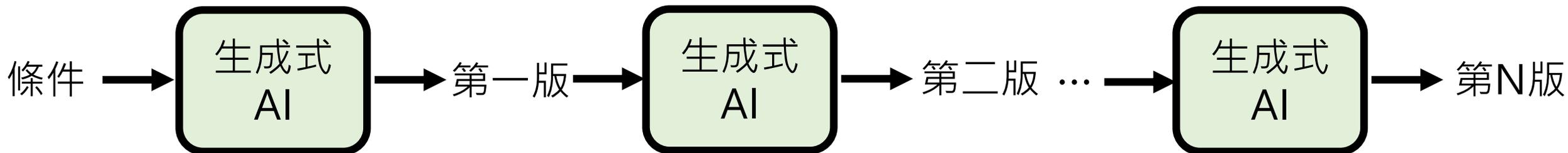
$$1440^2 \times (64 \times 64)$$



約 85 億次 attention

相差千倍

# 文字生影片的挑戰



第 K-1 版



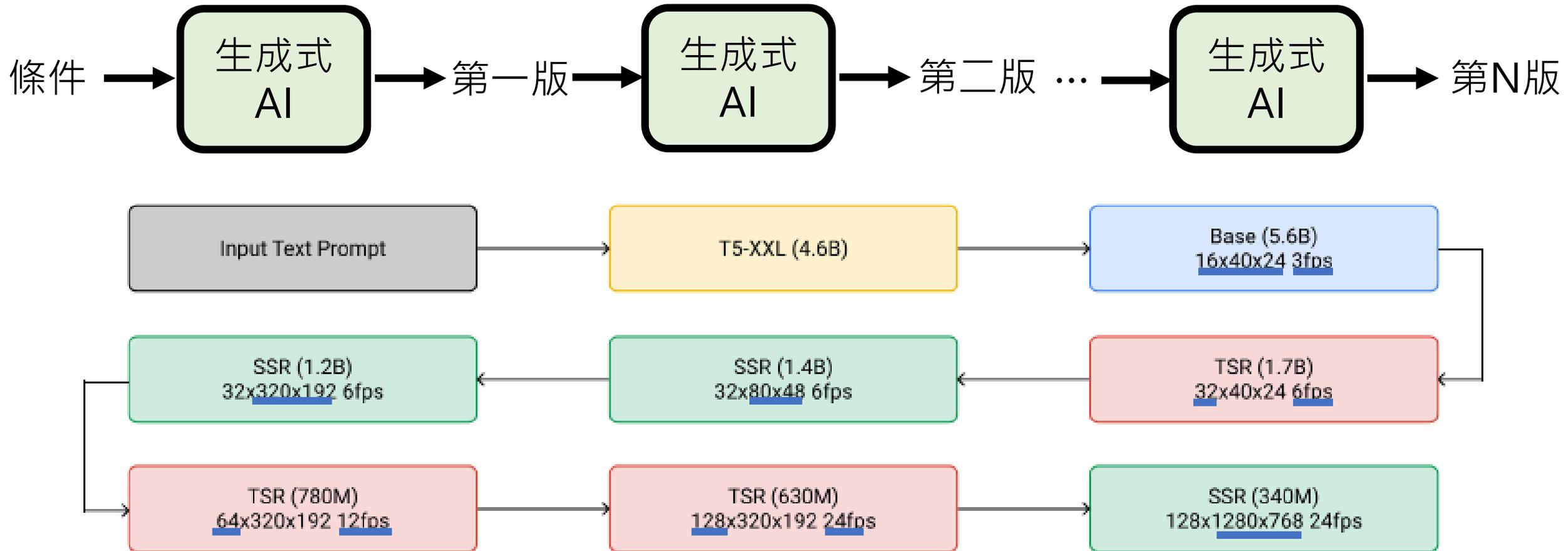
第 K 版

第 K-1 版



第 K 版

# 文字生影片的挑戰



# 延伸閱讀

<https://arxiv.org/abs/2405.03150>

## Video Diffusion Models: A Survey

**Andrew Melnik**

*Bielefeld University*

*andrew.melnik.papers@gmail.com*

**Michal Ljubljanac**

*Bielefeld University*

*mljubljanac@techfak.uni-bielefeld.de*

**Cong Lu**

*University of British Columbia*

*conglu@cs.ubc.ca*

**Qi Yan**

*University of British Columbia*

*qi.yan@ece.ubc.ca*

**Weiming Ren**

*University of Waterloo*

*w2ren@uwaterloo.ca*

**Helge Ritter**

*Bielefeld University*

*helge@techfak.uni-bielefeld.de*

# 經典影像生成方法介紹

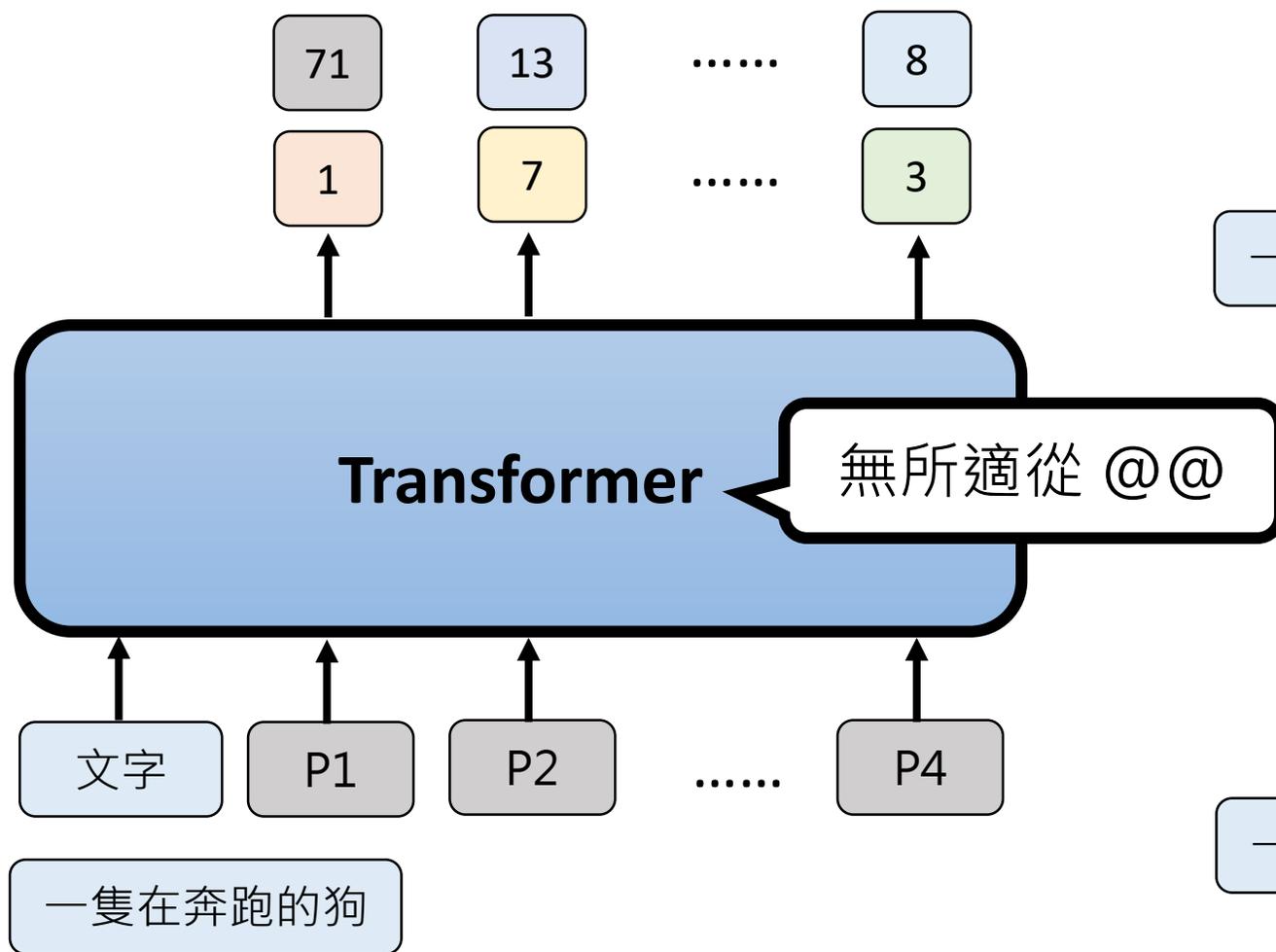
- Variational Auto-encoder (VAE)
- Flow-based Method
- Diffusion Method
- Generative Adversarial Network (GAN)

Sora 使用的是 Diffusion



Source of image: <https://openai.com/index/video-generation-models-as-world-simulators/>

# 文字生影像的挑戰



一隻在奔跑的狗



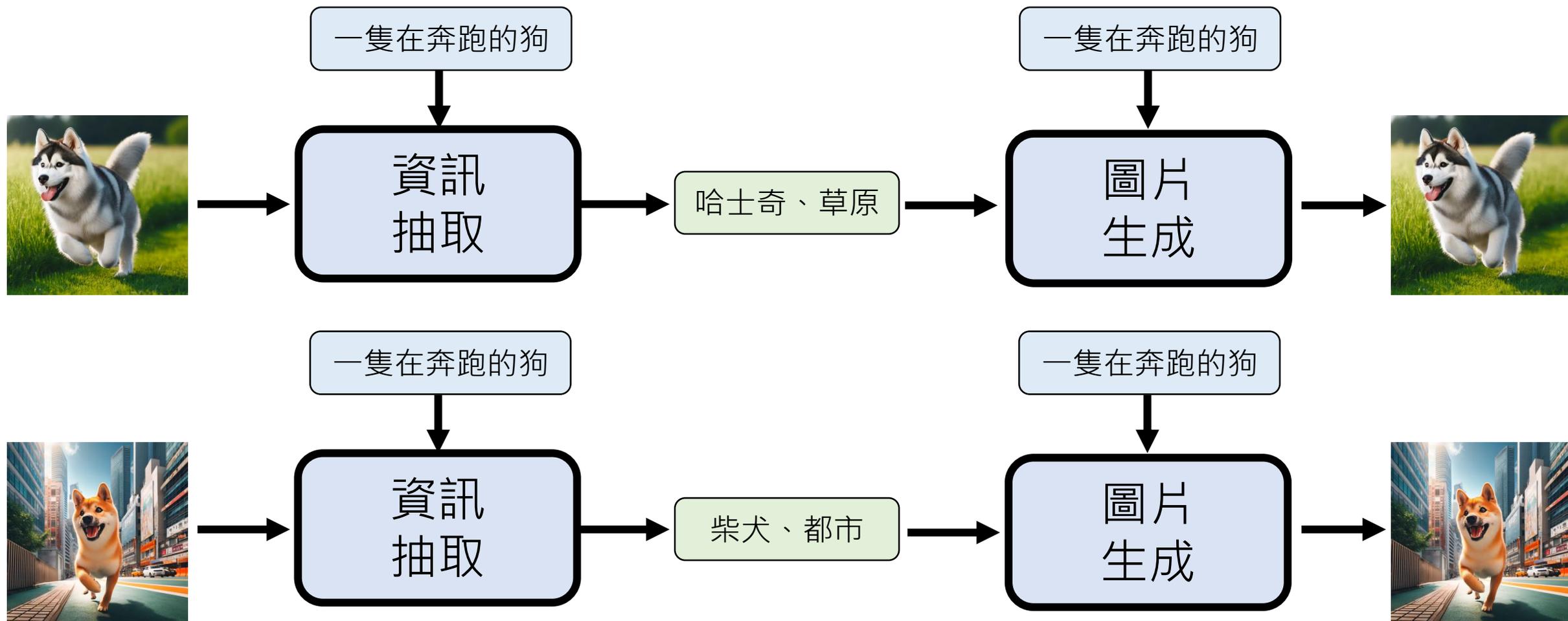
1 7 ..... 3

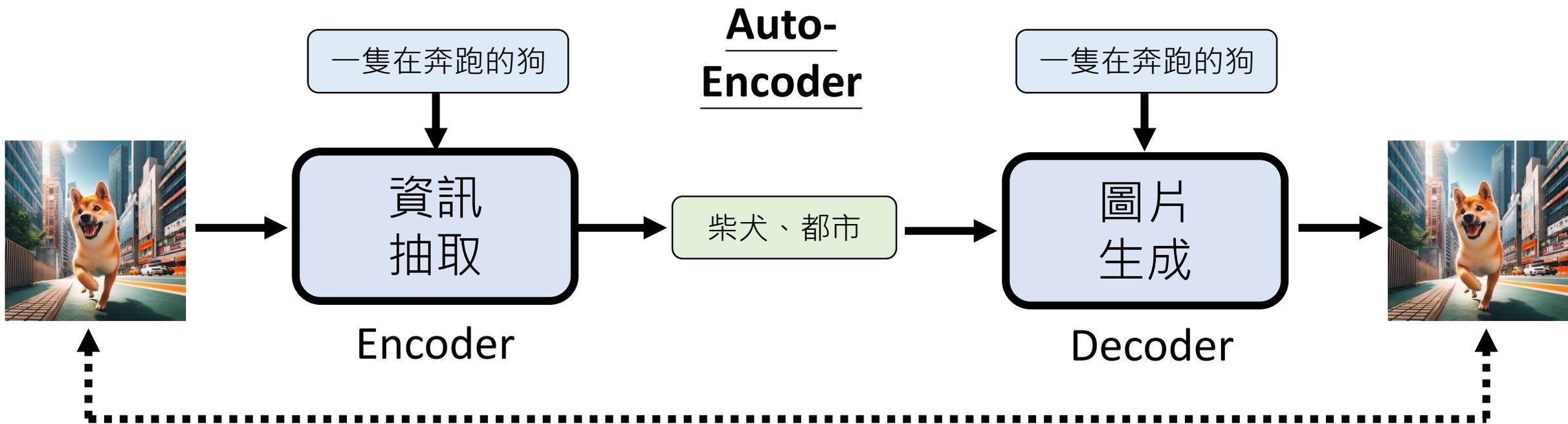
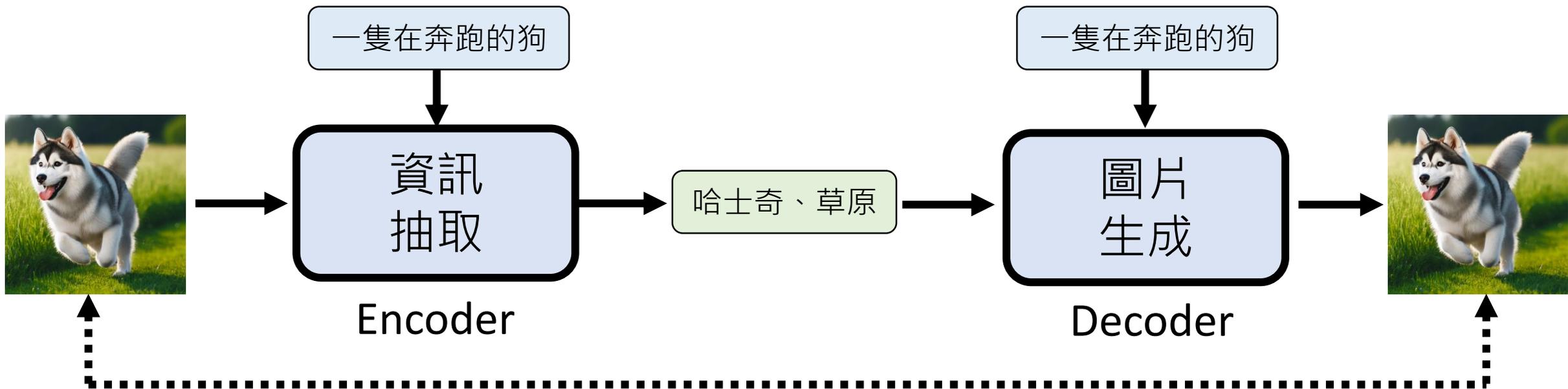
一隻在奔跑的狗

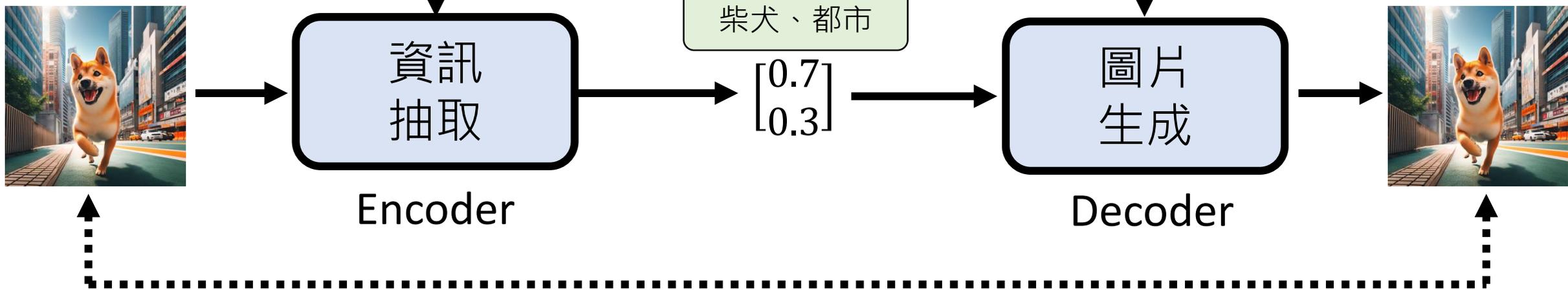
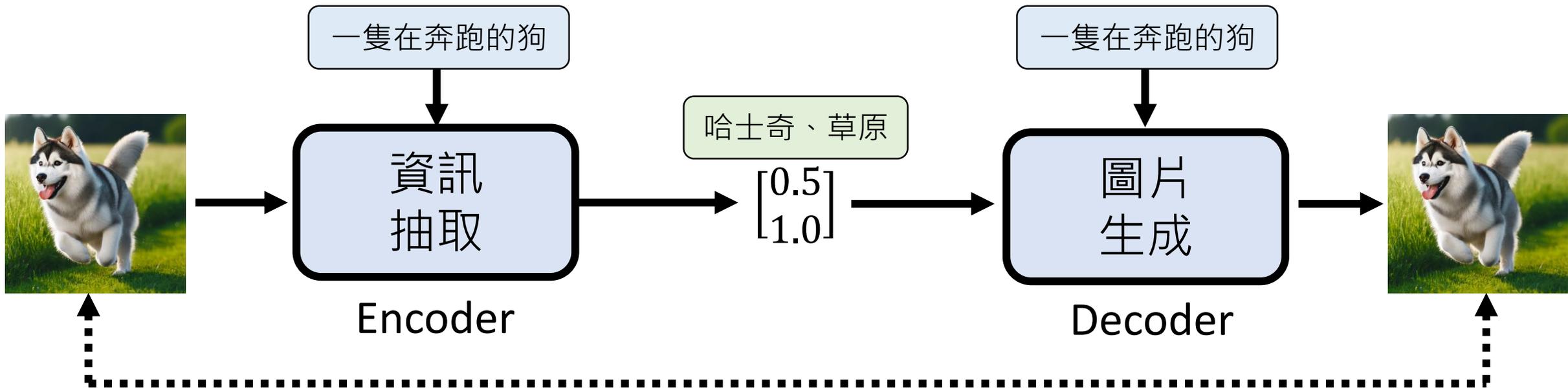


71 13 ..... 8

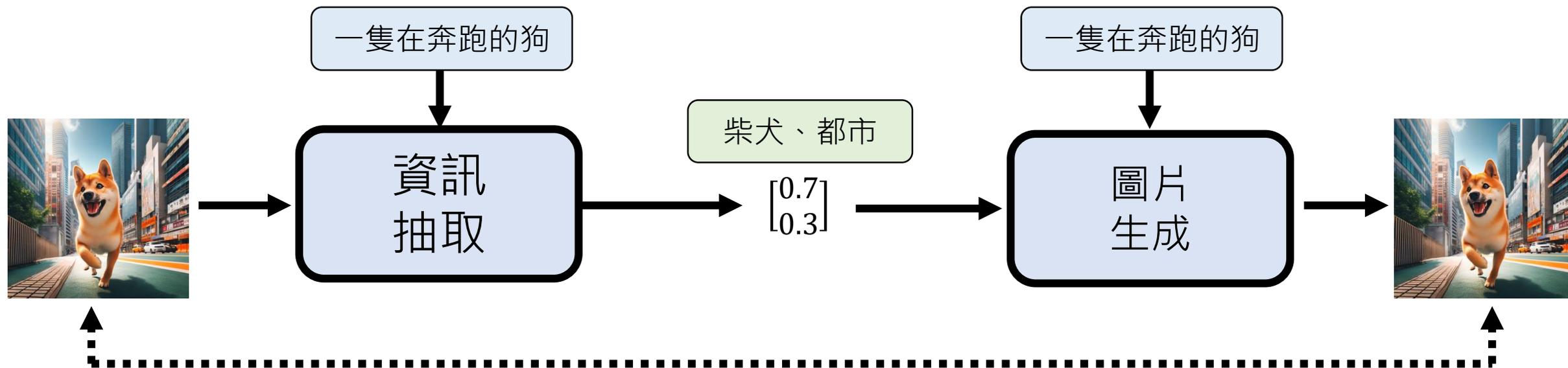
# 如何處理腦補問題



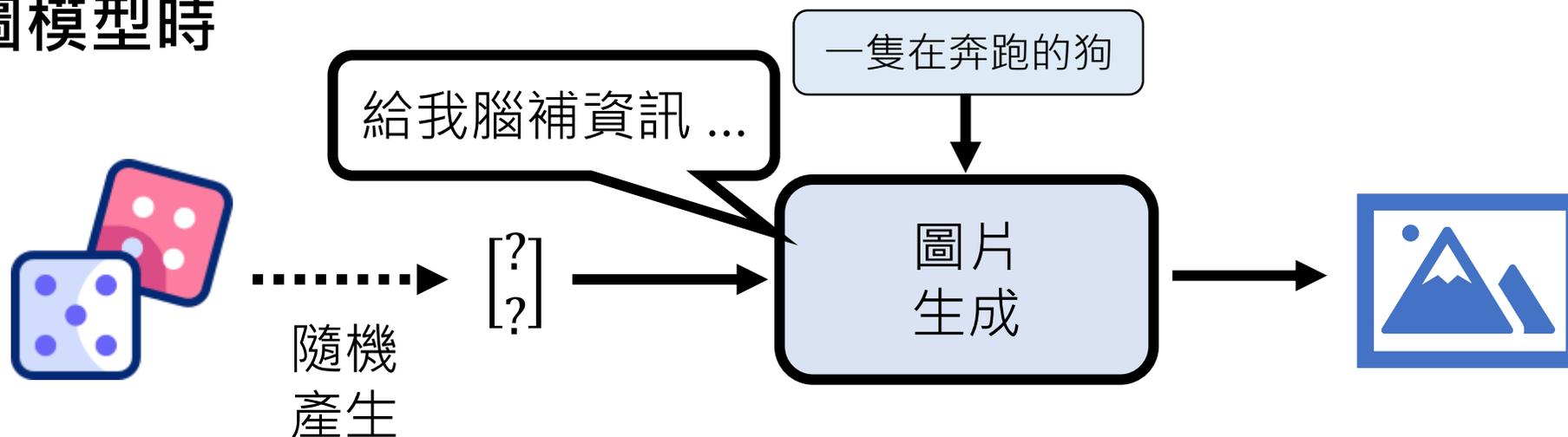




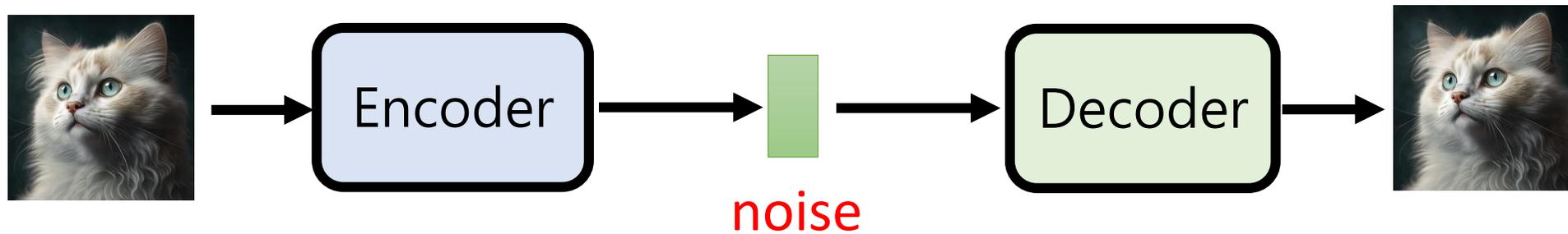
## 訓練文字生圖模型時



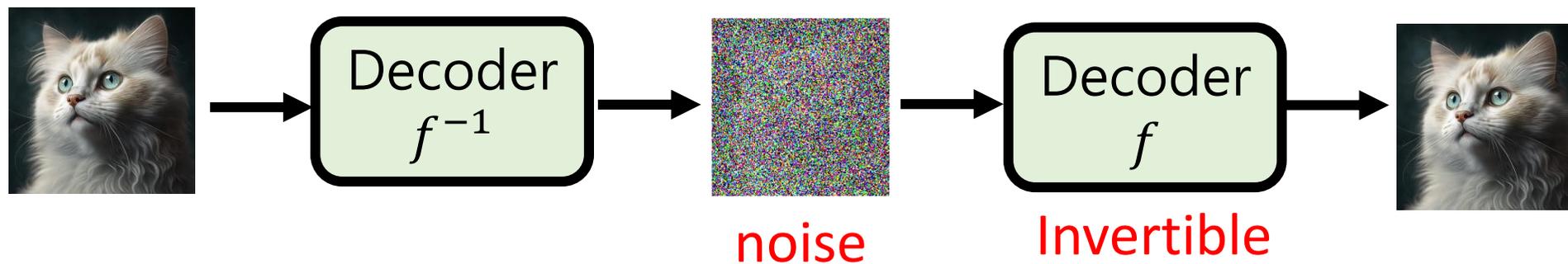
## 測試文字生圖模型時



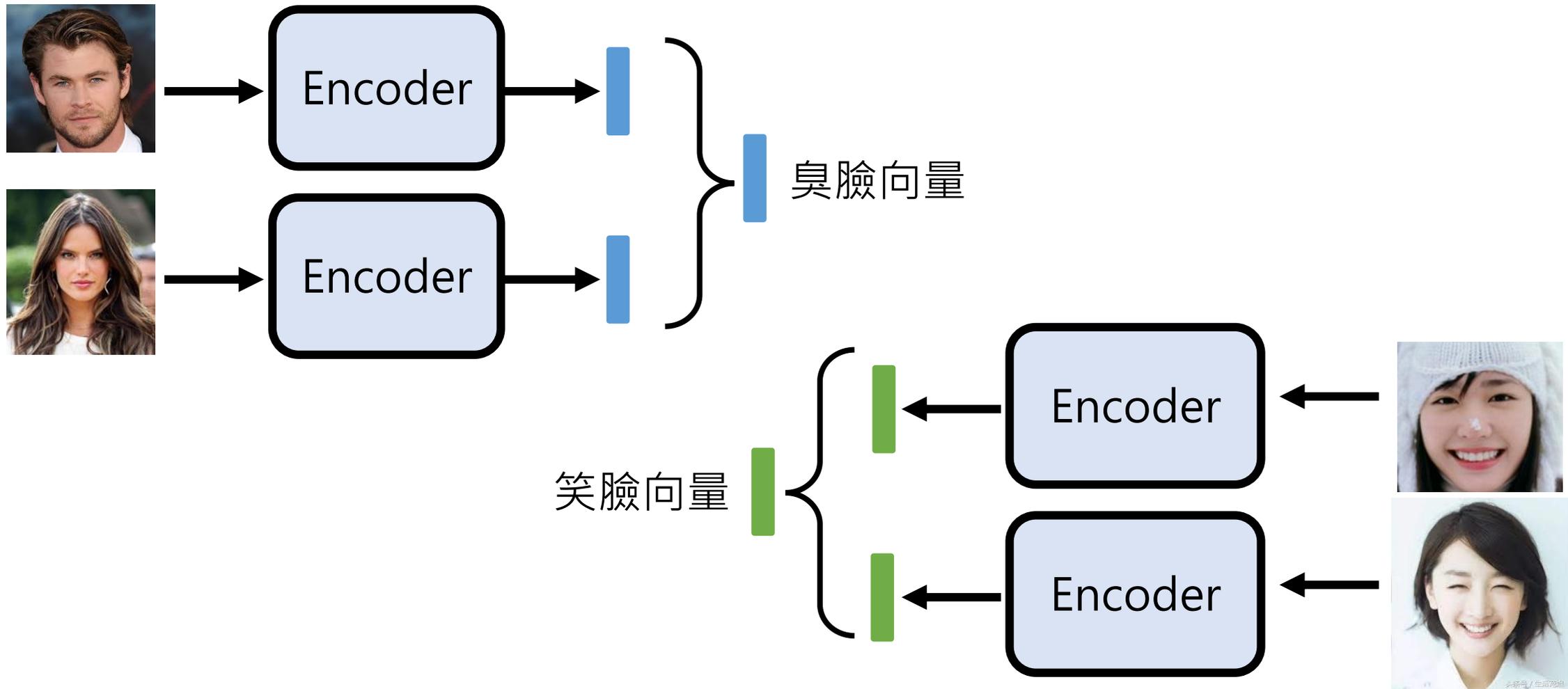
VAE



Flow-  
based

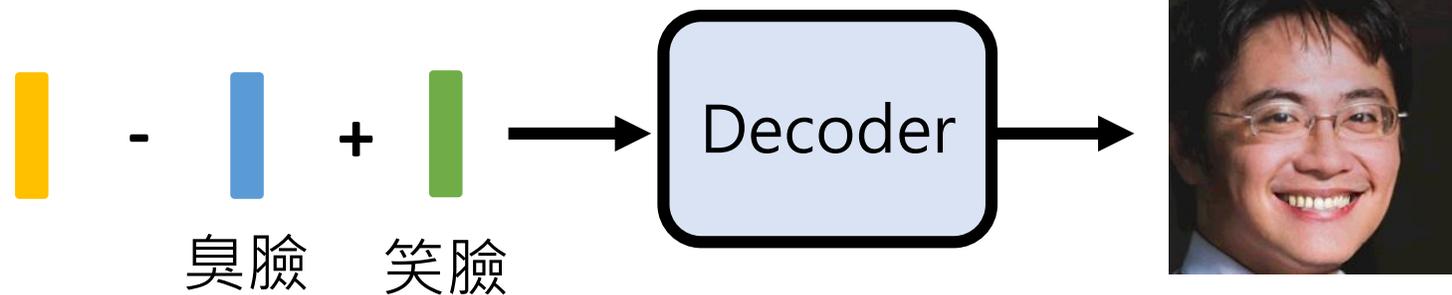
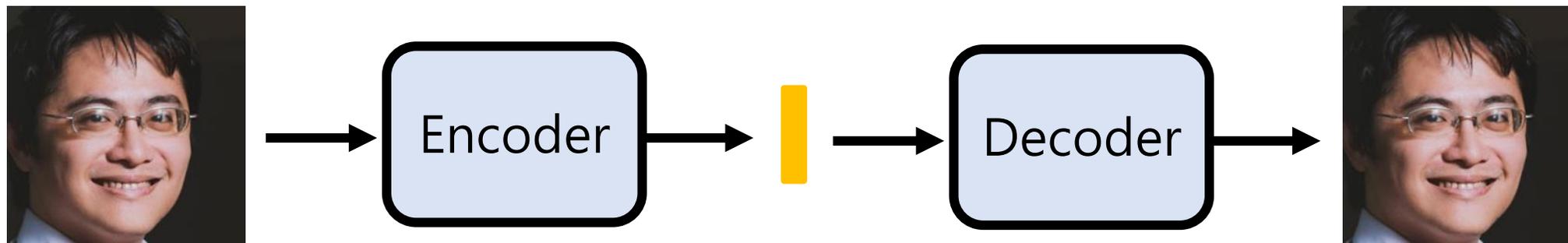


# 有資訊的雜訊 (noise)



# 有資訊的雜訊 (noise)

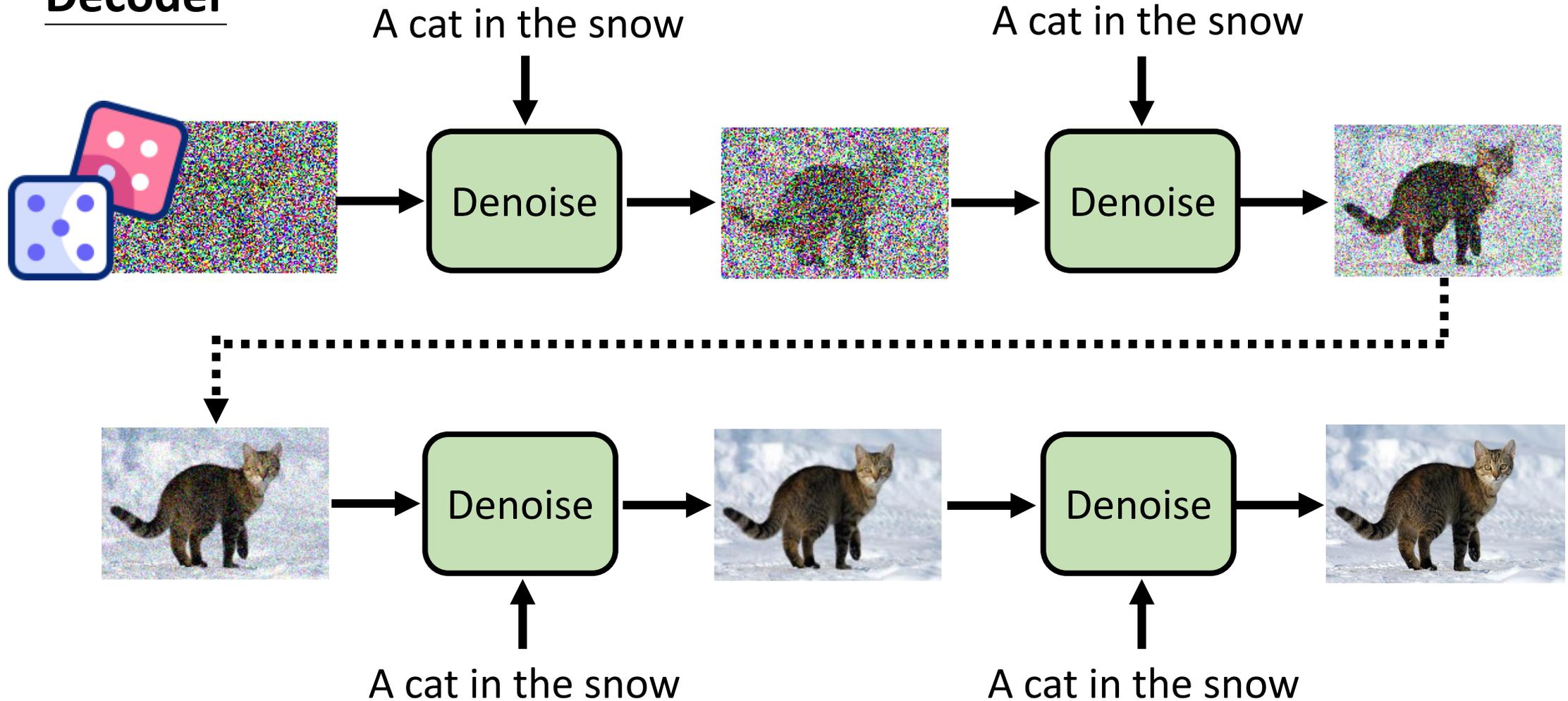
■ 臭臉向量    ■ 笑臉向量



此處是極度簡化後的講法，詳細說明請見：  
[https://www.youtube.com/watch?v=azBugJzmz-o&list=PLJV\\_el3uVTsNi7PgekEUFsyVIIAJXRSP-](https://www.youtube.com/watch?v=azBugJzmz-o&list=PLJV_el3uVTsNi7PgekEUFsyVIIAJXRSP-)

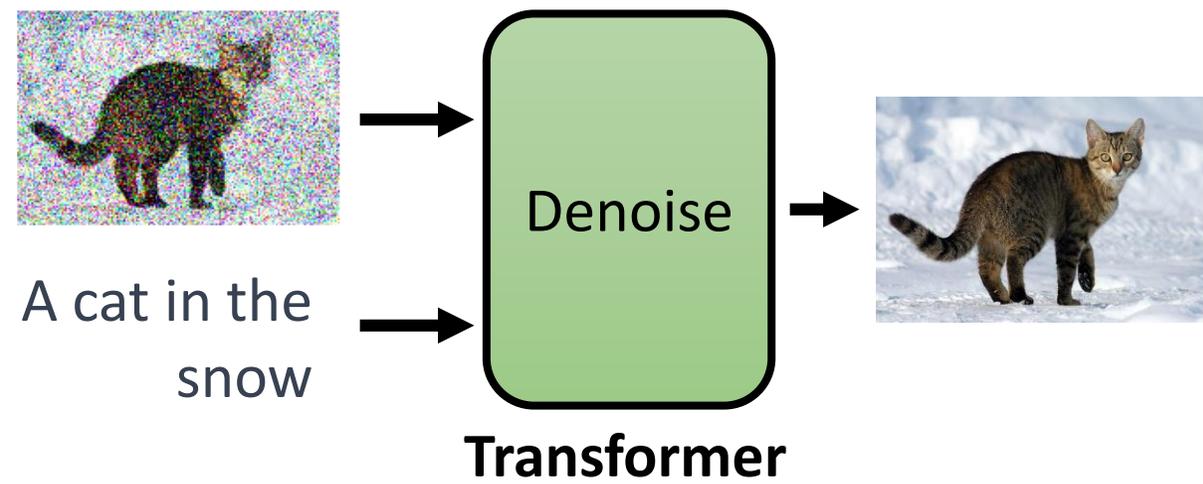
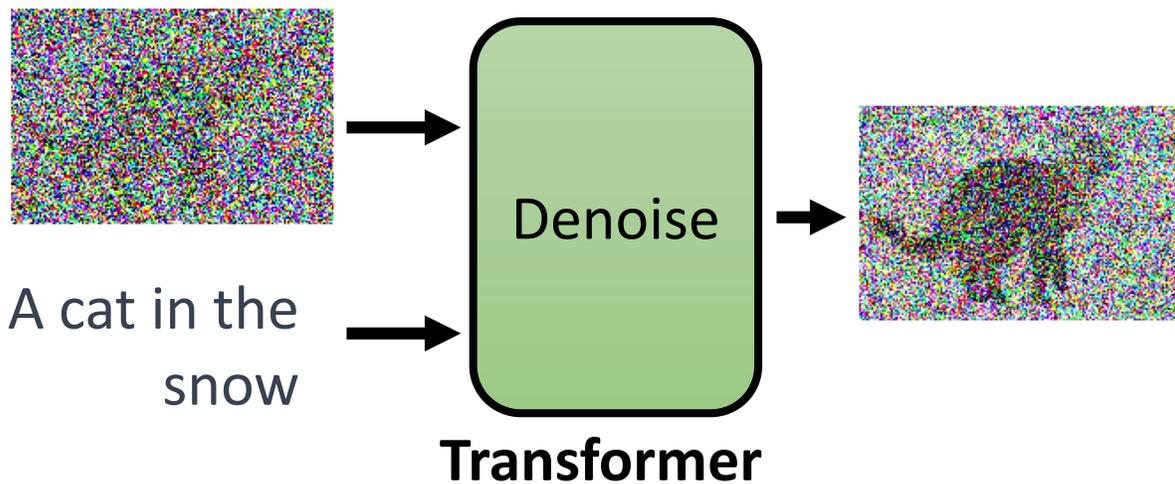
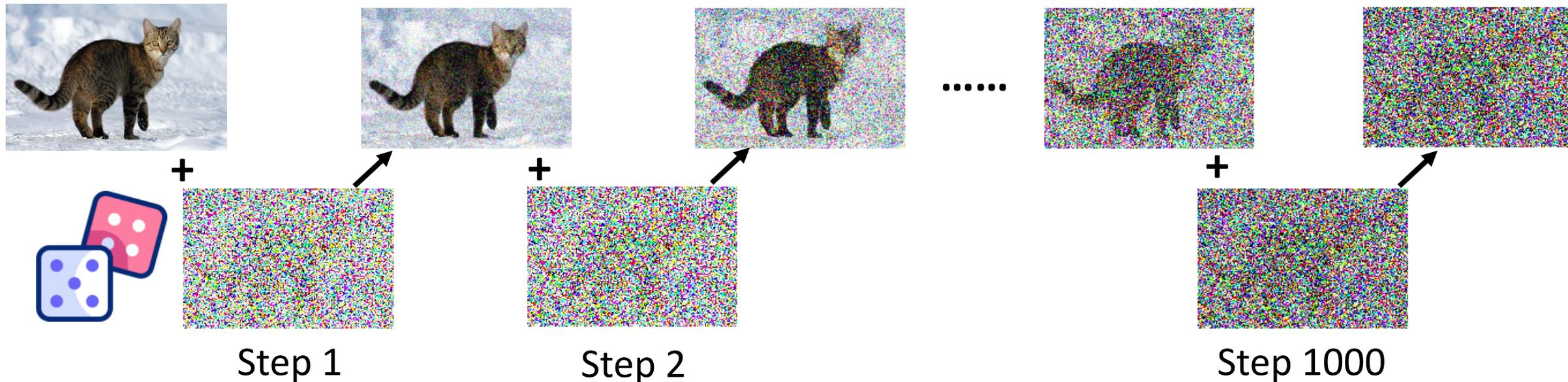
# Diffusion Model

## Decoder



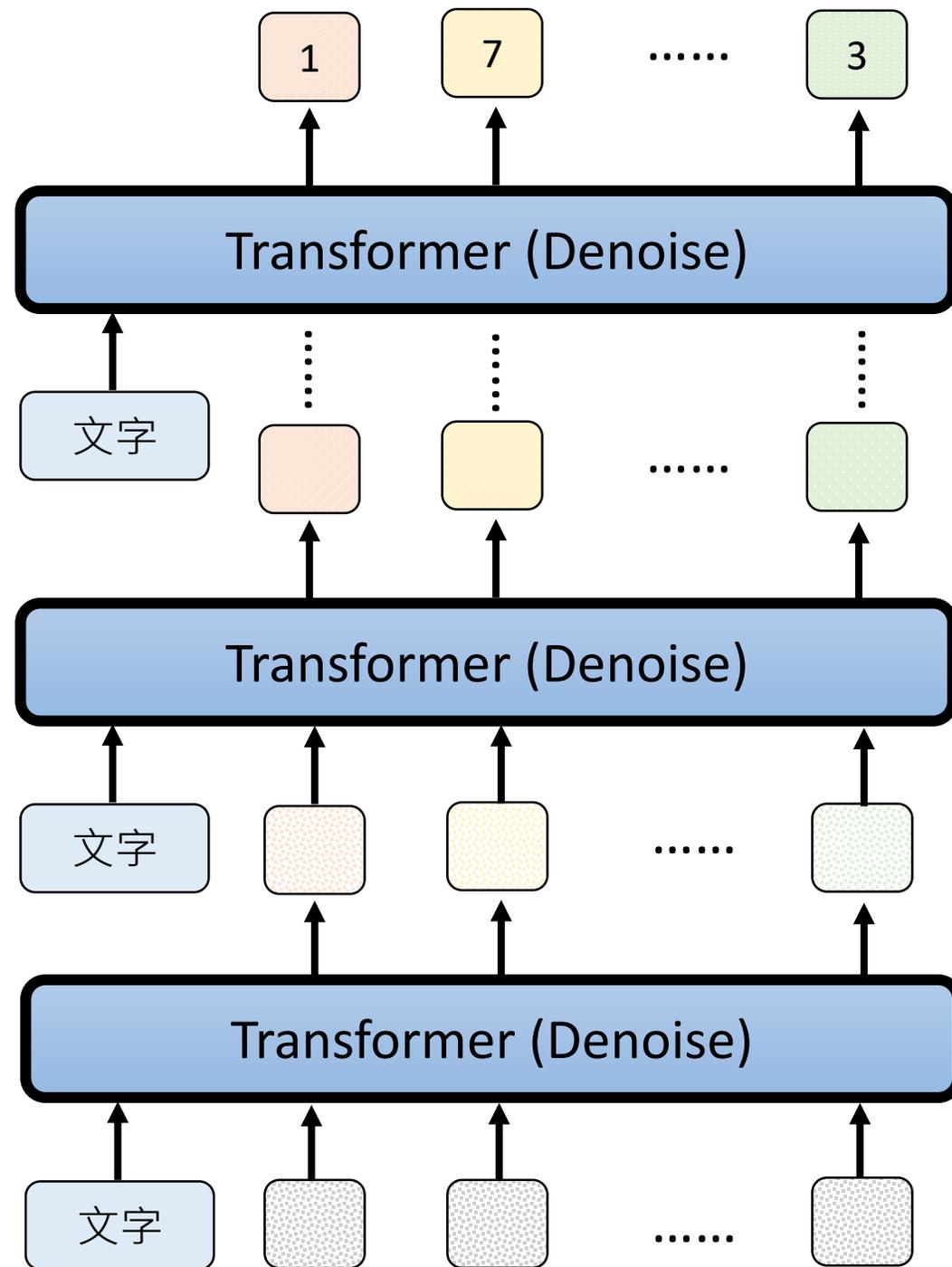
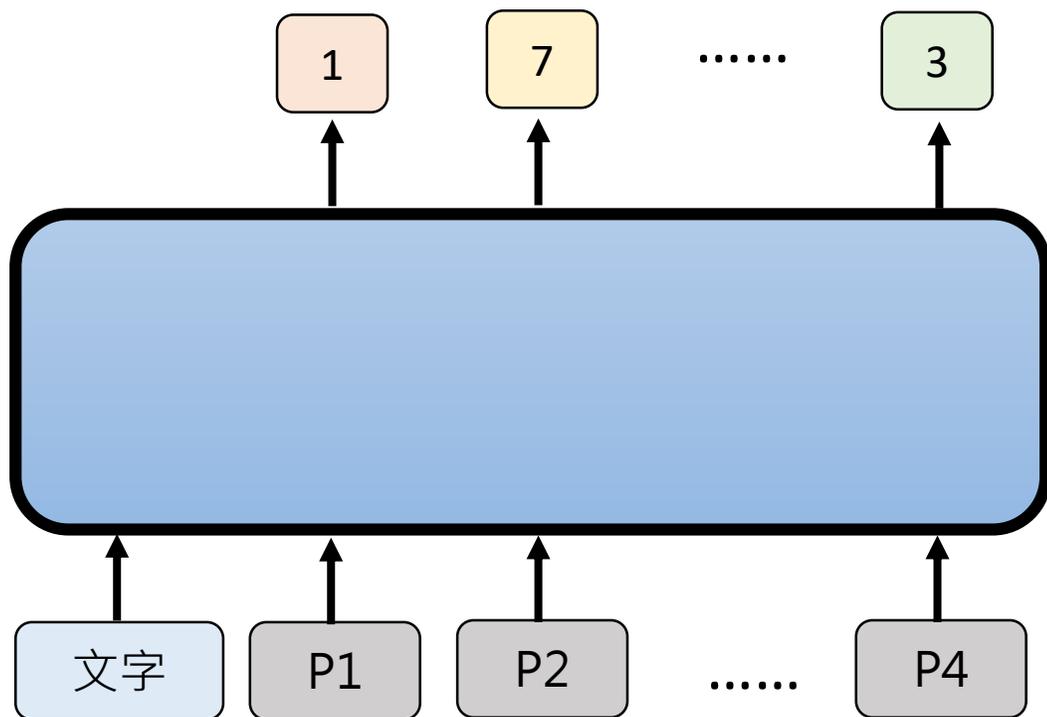
# A cat in the snow

此處是極度簡化後的講法，詳細說明請見：  
[https://www.youtube.com/watch?v=azBugJzmz-o&list=PLJV\\_el3uVTsNi7PgekEUFsyVIIAJXRSP-](https://www.youtube.com/watch?v=azBugJzmz-o&list=PLJV_el3uVTsNi7PgekEUFsyVIIAJXRSP-)



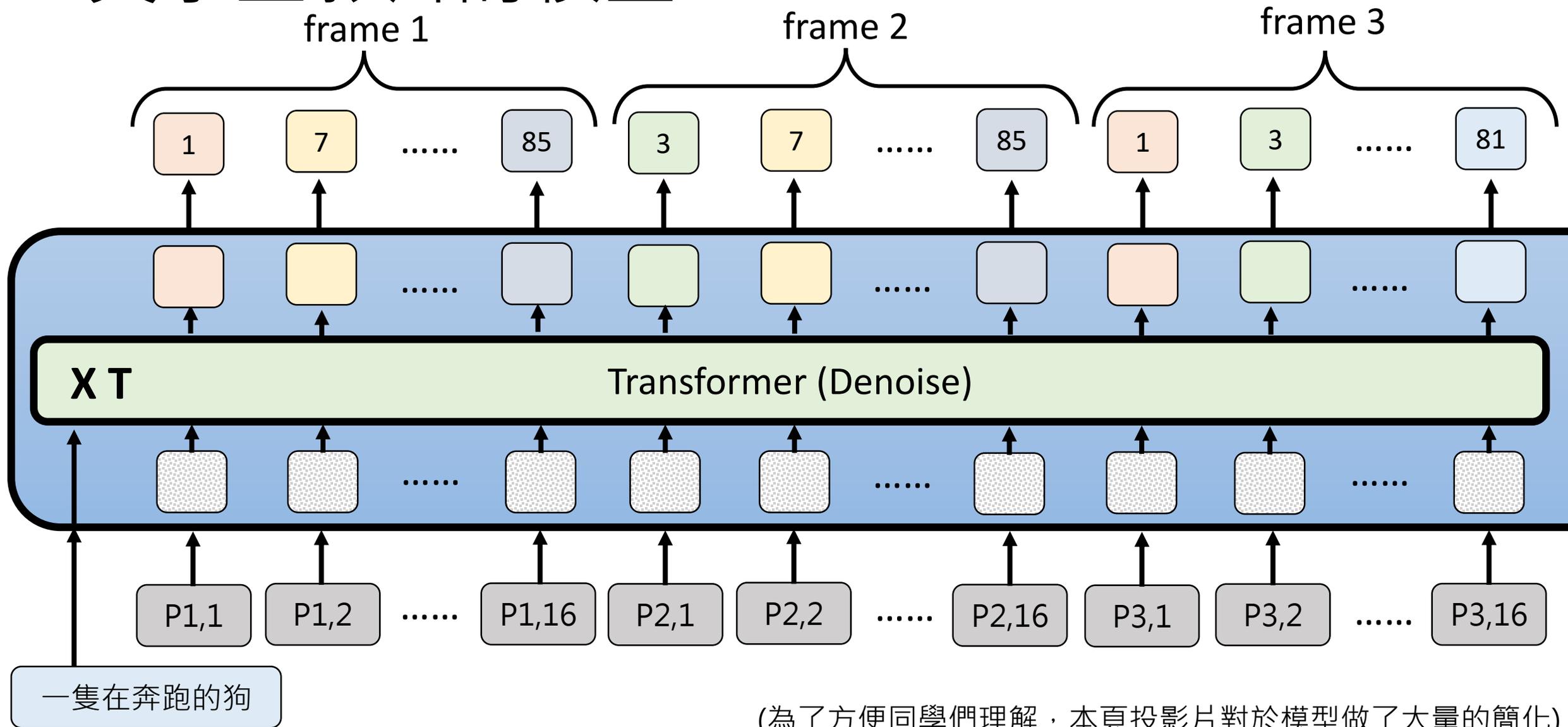
# Diffusion Transformer

<https://arxiv.org/abs/2212.09748>

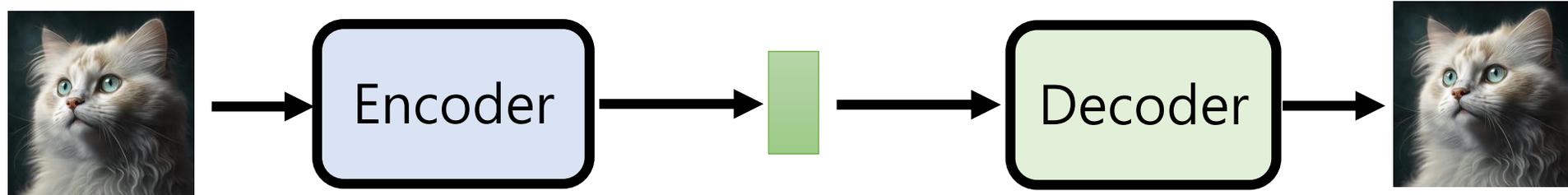


(為了方便同學們理解，本頁投影片對於模型做了大量的簡化)

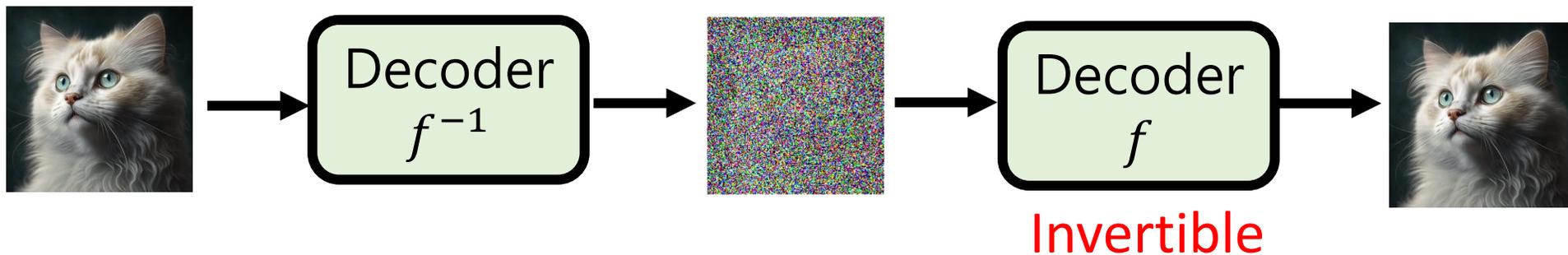
# 文字生影片的模型



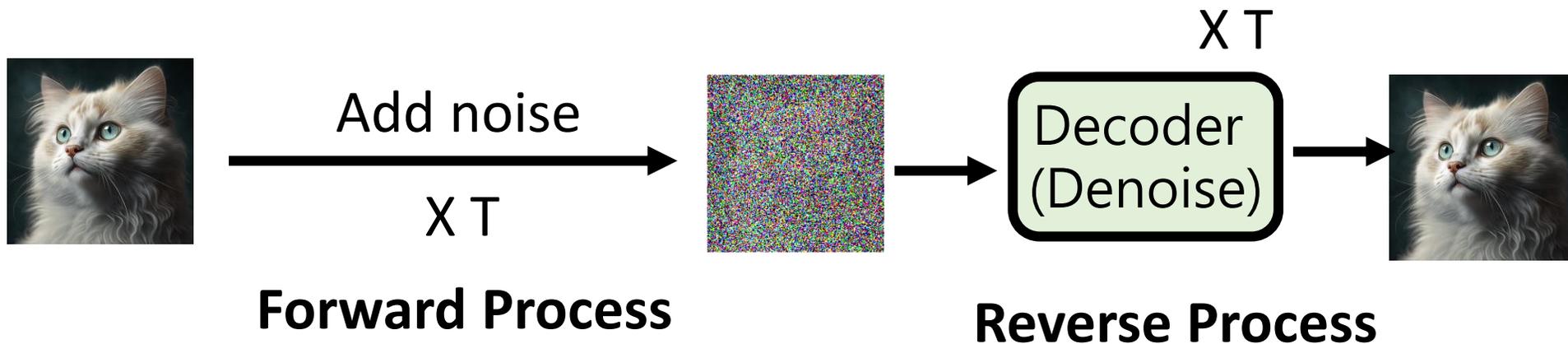
VAE



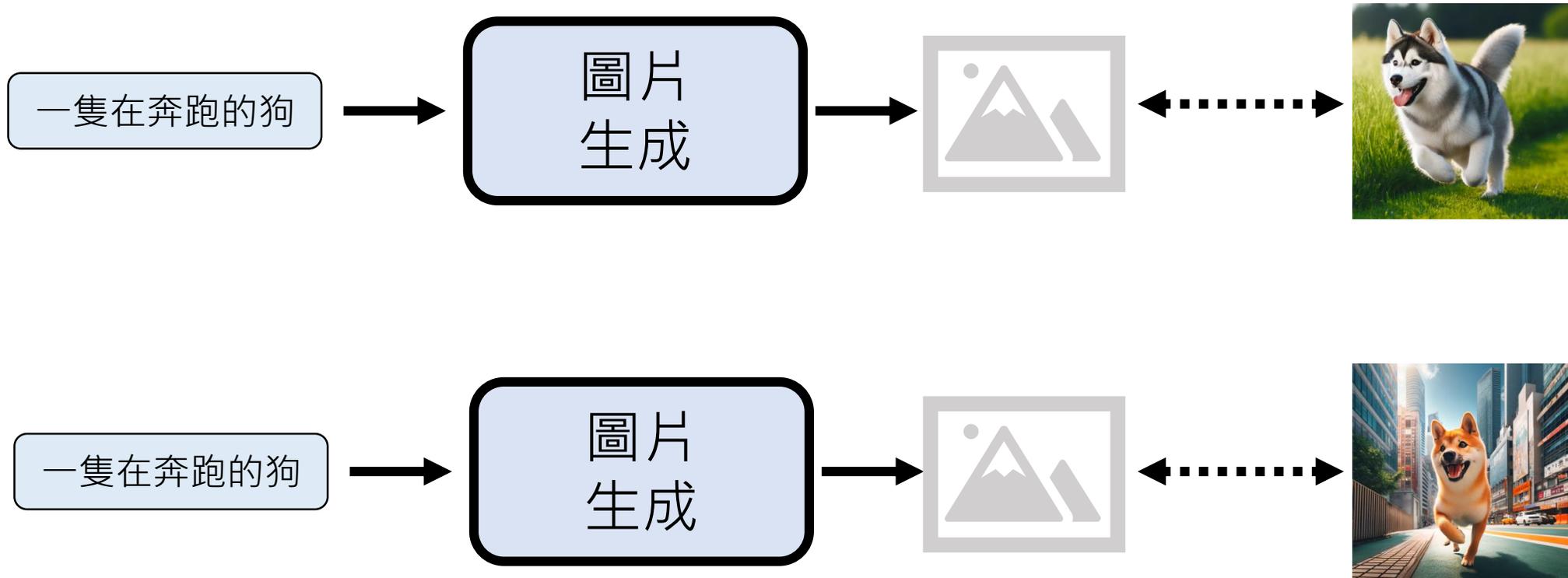
Flow-based



Diffusion



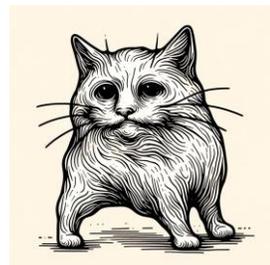
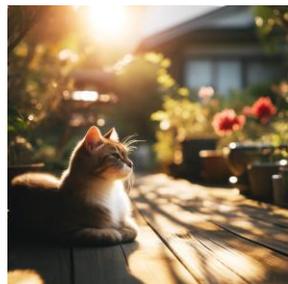
# Generative Adversarial Network (GAN)





一隻在奔跑的狗

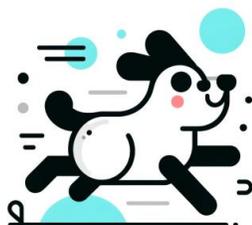
一隻在奔跑的狗



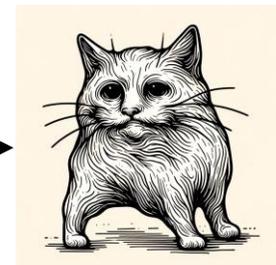
陽光下的貓

陽光下的貓

一隻在奔跑的狗

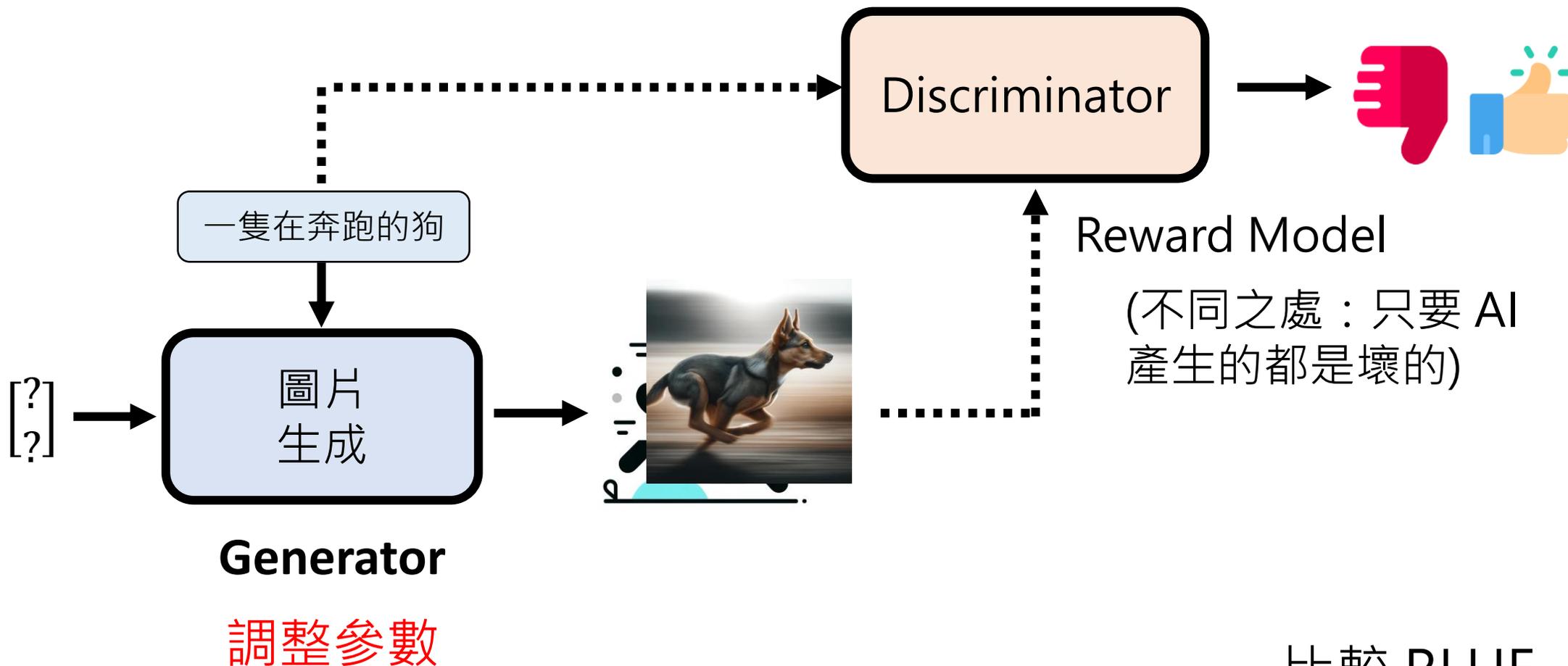


陽光下的貓



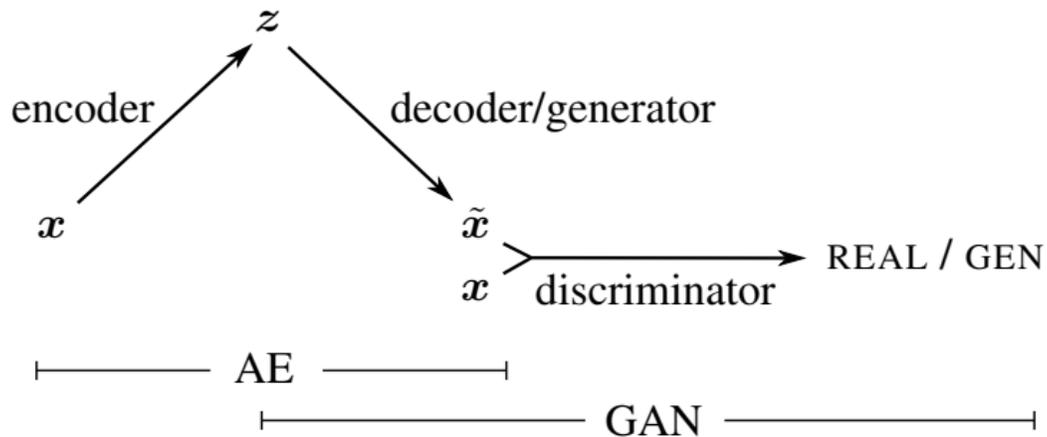
# Generative Adversarial Network (GAN)

Discriminator 和 Generator 會交替訓練



比較 RLHF

# GAN 是個外掛



## VAE + GAN

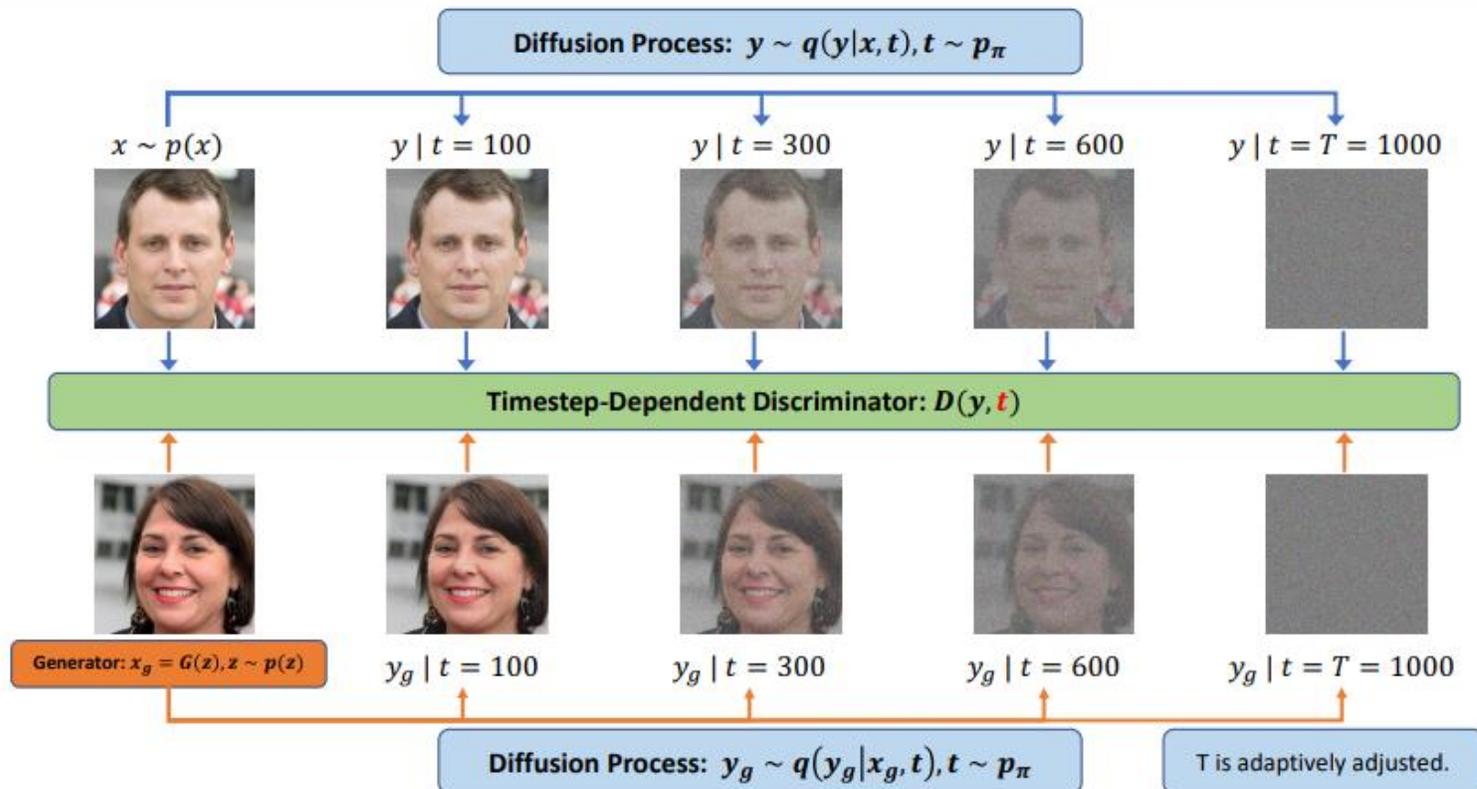
<https://arxiv.org/abs/1512.09300>

## Flow + GAN

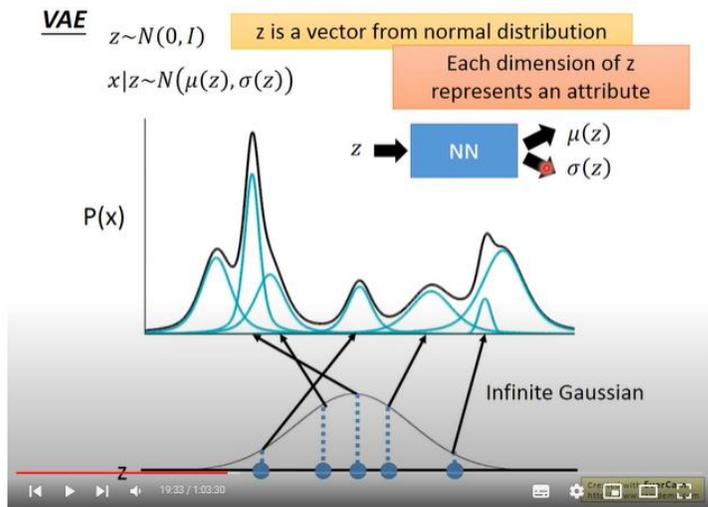
<https://arxiv.org/abs/1705.08868>

## Diffusion + GAN

<https://arxiv.org/abs/2206.02262>



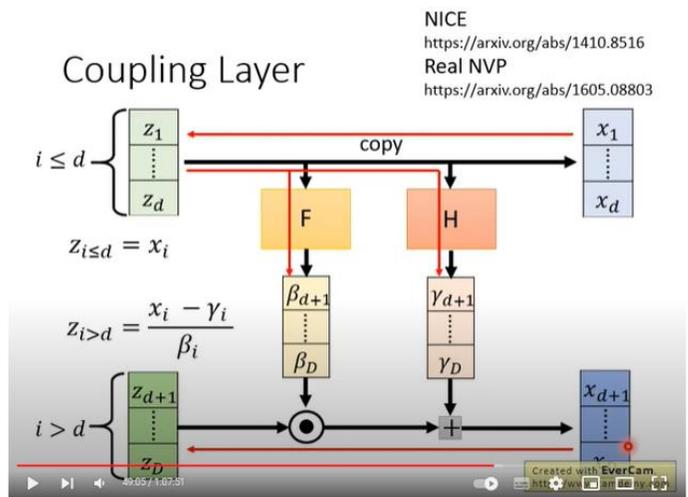
VAE



ML Lecture 18: Unsupervised Learning - Deep Generative Model (Part II)

<https://youtu.be/8zomhgKrsMQ>  
(2016 機器學習)

Flow-based



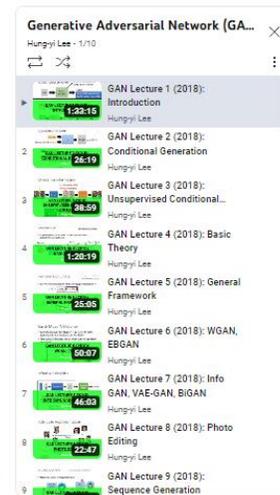
Flow-based Generative Model

<https://youtu.be/uXY18nzdSsM>  
(2019 機器學習)

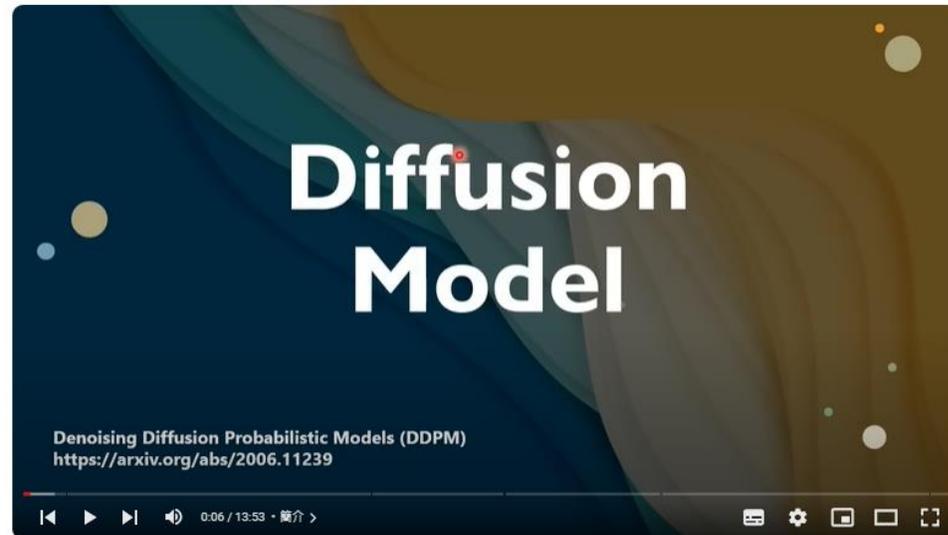
GAN

# Introduction of Generative Adversarial Network (GAN)

李宏毅  
Hung-yi Lee



[https://www.youtube.com/watch?v=DQNNMiAP5lw&list=PLJV\\_el3uVTsMq6JEFPW35BCiOQTsoqWnW](https://www.youtube.com/watch?v=DQNNMiAP5lw&list=PLJV_el3uVTsMq6JEFPW35BCiOQTsoqWnW) (2018 機器學習及其深層與結構化)



Diffusion

[https://www.youtube.com/watch?v=azBugJzmoz-o&list=PLJV\\_el3uVTsNi7PgekEUFsyVIIAJXRSP-](https://www.youtube.com/watch?v=azBugJzmoz-o&list=PLJV_el3uVTsNi7PgekEUFsyVIIAJXRSP-) (2023 機器學習)



# 有沒有可能跟生成的影像有更強的互動

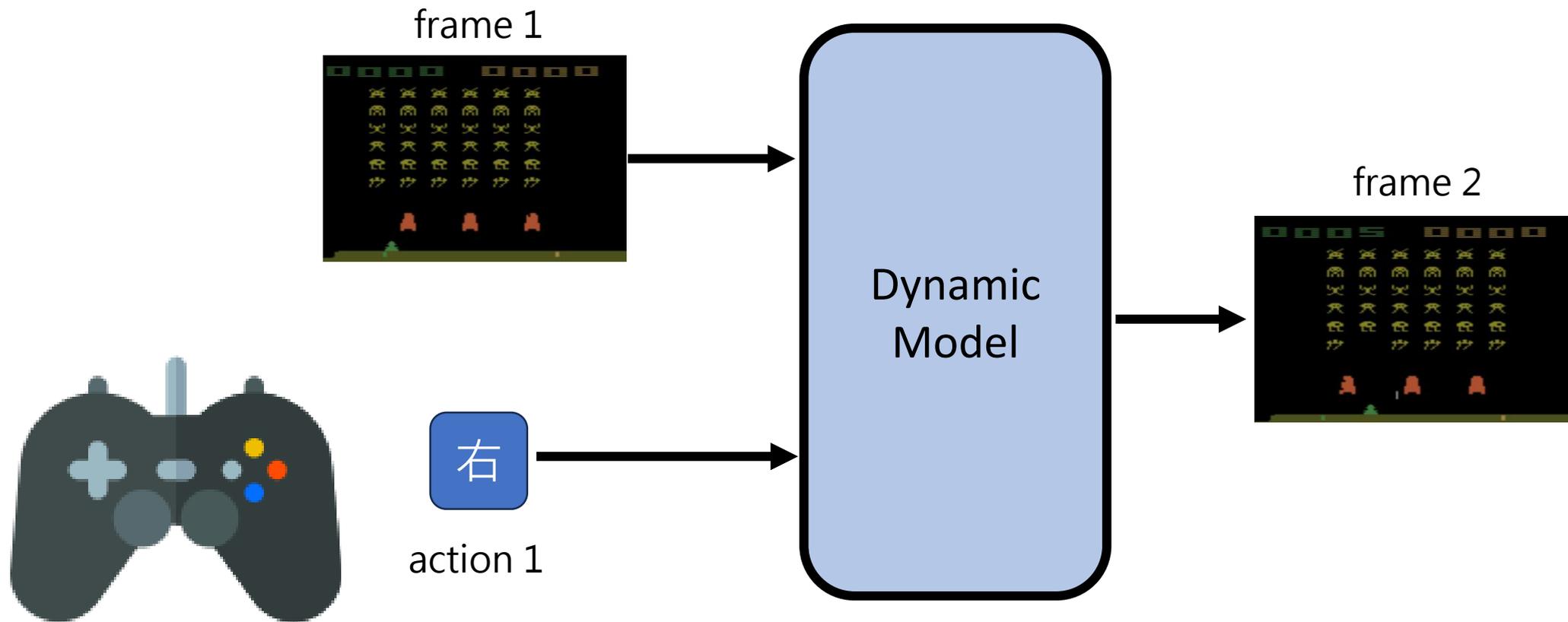


例如：直接操控這個人要走去哪裡

直接做個開放世界遊戲？

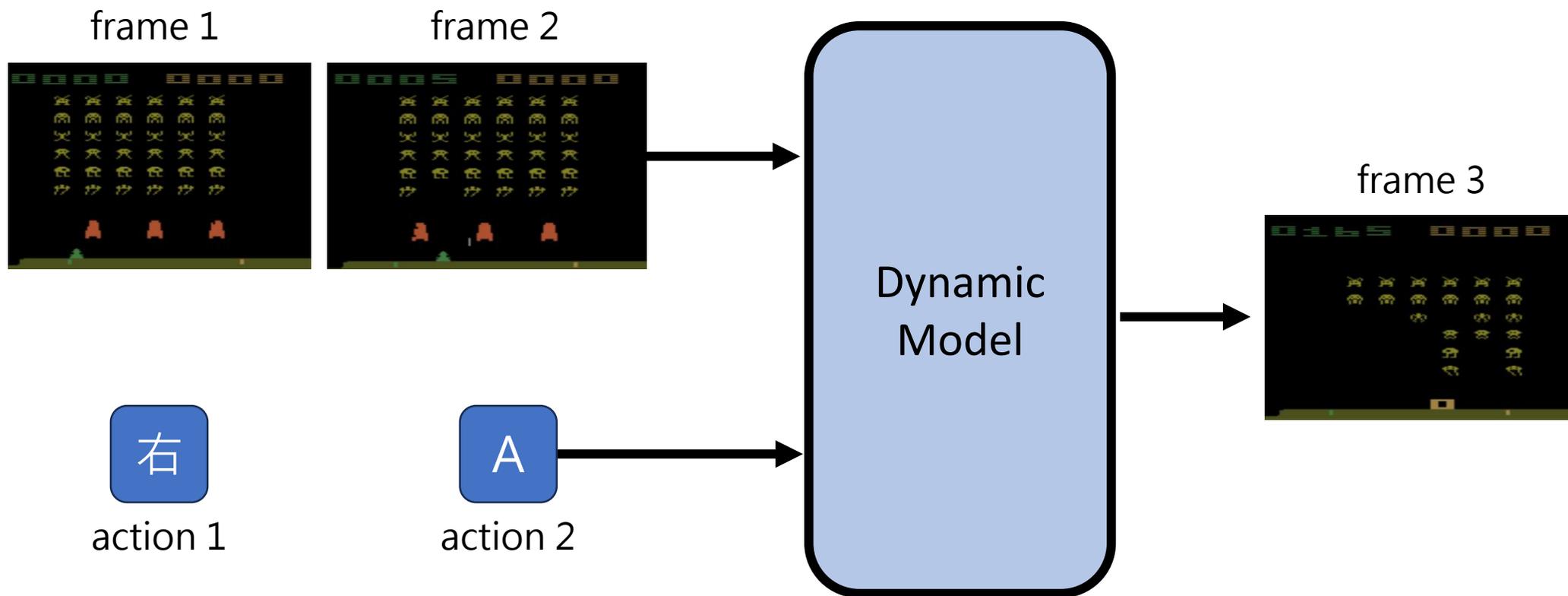
# Genie: Generative Interactive Environments

<https://arxiv.org/abs/2402.15391>



# Genie: Generative Interactive Environments

<https://arxiv.org/abs/2402.15391>



# Genie: Generative Interactive Environments

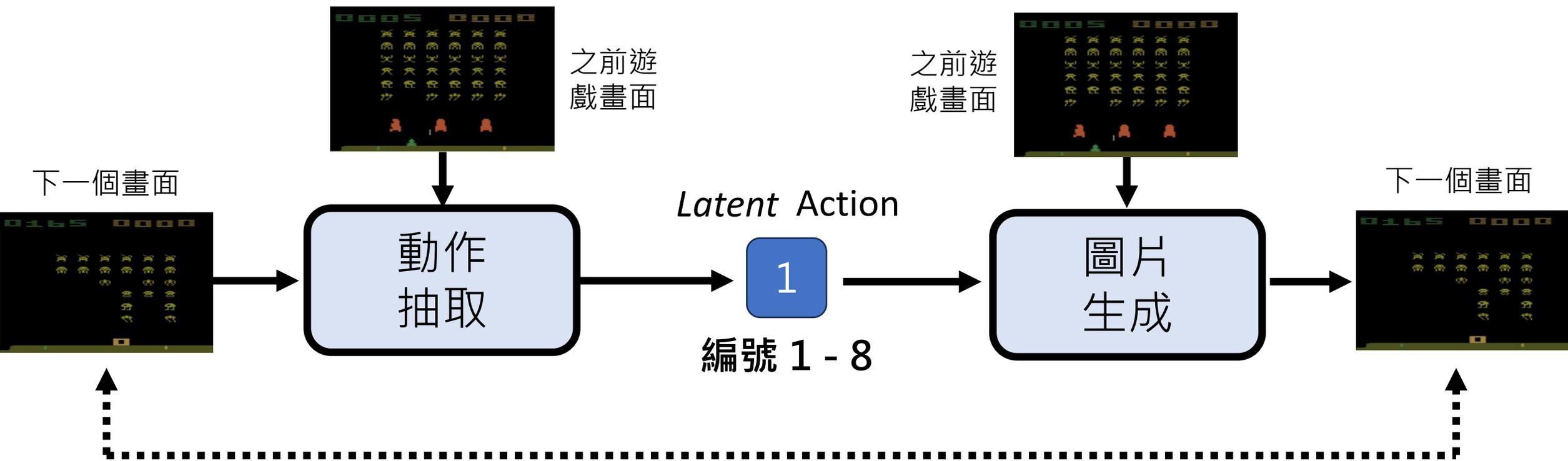
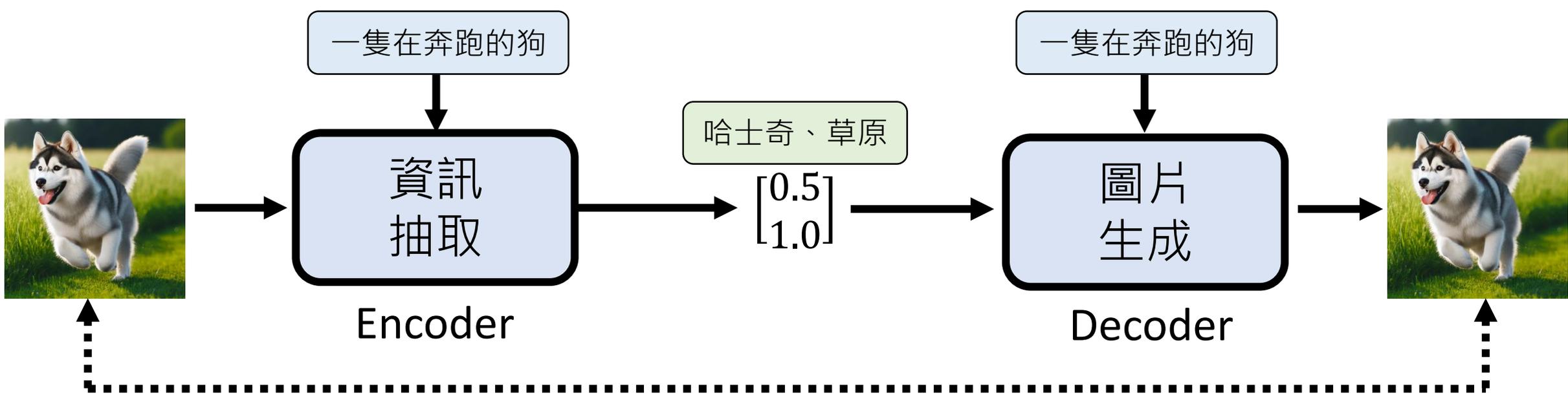
- 訓練資料？



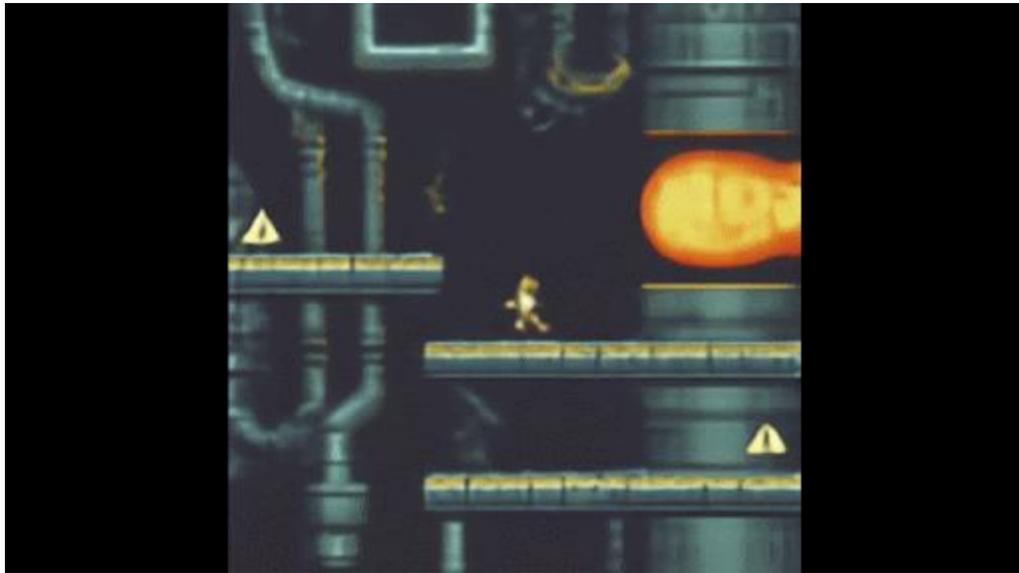
可以蒐集大量  
遊戲影片



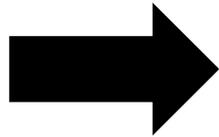
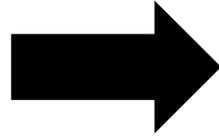
沒有使用者輸入的動作



# Genie: Generative Interactive Environments



Latent action: 6, 6, 7, 6, 7, 6, 5, 5, 2, 7



Source: <https://sites.google.com/view/genie-2024/home>

Sora 技術  
報告標題

February 15, 2024

# Video generation models as world simulators

[View Sora overview](#)

<https://arxiv.org/abs/2309.17080>

<https://wayve.ai/thinking/scaling-gaia-1/>