# Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
   A. **이름: 신준원**
   B. **저널: CVPR**
   C. **도메인: Diffusion**
   D. **출판연도: 2023**
   E. **저자: Jing Nathan Yan, Jiatao Gu, Alexander M. Rush**

2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. **논문 Figure를 그대로 따라 그리면 안됩니다.**
   A. [선행연구] Attention Block은 U-Net, Transformer를 사용할 때 빠짐없이 사용됨. 그러나, 해당 Block은 고해상도에서 높은 연산량, 파라미터 수를 요구한다는 단점이 있음.

   B. 따라서, 해당 논문에서는 Attention Block 없이 학습가능한 Architecture를 제안하고 있음.

   C. SSM (State Space Model) Block을 제안하고 있음. Transformer의 경우 Patchify를 요구하고 있는데, 이 과정에서 High-frequency 표현, Patch 결합 과정의 불안정성 등의 문제가 있다고 판단하고 있음. 따라서, Patch화 하지 않은 상태에서, Attention Block을 제거할 수 있다는 점에서 SSM Block은 우수한 선택이라 할 수 있음. RNN 형태의 구조와 닮았으나, SSM의 Equation에 따라, 연산량을 감소시킬 수 있음(sub-quadratic)과 동시에, Convolutional한 계산이 가능함.

   D. 또한, Input Image(Noised)를 직접 사용하지 않고, Latent Diffusion의 Encoder를 사용, 이를 Flatten화 해서 사용함(DiT와 동일- Patch, MLP).

   E. Condition(C, T)는 Scale-Shift 형태로 적용하며, Attention은 Elemental-wise sum + Mul 형태로 처리하고 있음.

   F. 결과적으로, Linear Decoder를 사용해, Input 과 동일한 Size로 통일해줌(DiT

와 동일 -> Noise, Covariance Prediction).

3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet을 그대로 따라 적으면 안됩니다.**
   A. Attention Block을 제거했음에도, 비슷한 성능을 뽑았다는 점은 유의미하다고 판단됨.

4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
   A. 고해상도 문제 Attention의 연산량이 높다는 점을 문제로, 마치, U-Net을 Transformer로 대체했던 것처럼, 이를 효과적으로 대처했다는 점에서, 현재 정답이라고 여겨지는 구조에 한 번씩 의문을 가져보는 시도가 필요해 보임.

5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
   A. Diffusion Transformer, DiT-XL 대비 연산량은 증가했으나, 성능적으로는 부족한 모습을 보인다는 점이 아쉬움.

6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
   A. Diffusion 선행연구

7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
   A. 아니요

날짜: 2025-07-16

이름: 신준원

## introduction

Diffusion → 연산량이 너무 높다. ⇒ Diffusion state space model (DiffuSSM)

연산량 많은 Attention ⇒ State space model backbone으로 대체

↓

U-Net, transformer 모두 필요.

① representation Compression                    Spatial information

→ patchifying.    ⊞    단점 → high-frequency ↓

(trade-off) structural integrity ↓

나눠서 정보 ⇒ 원본대비 파악 어려워짐.

→ multi-scale resolution | 단점 → Spatial detail ↓ (down-sample)

attention 연산량↓.          o Generate. (격자) (up-sample)

~~Channel~~ → 저장영향                    resolution

⊢ sub-quadratic space

(one of Attention Approximation method)

② DiffuSSM              ⊢ (LRA. Audio → 사용됨)

→ gated state space model (SSM) backbone 활용

↓

Sequence model 성능 향상 기여 효과.

⊕ hourglass architecture

(sequence 길이 처리)          VAE (latent model)의 encoder 사용.

▷ z² (representation) ≒ flatten 하여 처리

# State space models (SSM)

input sequence of scalars = $u_1, u_2, u_3, \ldots, u_L$ / output = $y_1, y_2, \ldots, y_L$
($L \times 1$)                                         $1 \times N$.

equation $\Rightarrow$ $d_k = \overline{A} d_{k-1} + \overline{B} u_k$ , $y_k = \overline{C} d_k$

               $N \times N$    $N \times N$.    $N \times 1$      $1 \times N$   $1 \times N$.

장점 : long convolutional 연산 가능 (with linear equation)

특징 : ① FFT (Fast Fourier Transform)

       └▷ sub-quadratic Algorithm.

            └▷ $O(N^2 b)$ ⇒ 해당 Algorithm $O(L \log L)$     } ▷ long-sequence handling에 용이함

     ② $\overline{A}, \overline{B}, \overline{C}$ ( discrete-time values )

         └▷ Continuous-time state-space 에서 가느라면.

       ⇒ stable & effective한 결근 가능

따라서, 논문에서는 ⇒ A,B,C 를 Continuous한 상태에서 discrete 하게 변경하는 것 사용 (S4D)

      ⊕ Bidirectional SSM layer 사용

        ⇒ RNN ~ Bidirection.

            └▷ Flatten, Global feature 반영 & 문제 해결를 위한 사용.

**Diffussm Block**

구성요소

① Gated bidirectional SSM  ⟶ flatten input
② hourglass Architecture (with ⓜⓛⓟ)  ⟶ m/p기반.
⟶ Sequence.

Block소개  V2차 방식  ⟶ expanding ⟷ Contracting.

I  ⟶ Hour glass

$J \times D$  ⟶ 긴해성도 (long length)
⟶ M = L/J
변해용  ⟶ 짧해성용 (low length)

C, t  flatten. ⟶ m/p based 방식 계산.

Condition

noise input  ⟶ Encoder ( VAE — Latent diffusion)

□ □ - - - □ □ □  ( z - sequence )
⟶ flattened.

Diffussm Block

Shift Scale-down | dense  ⟶  Condition → Scale-shift → Hourglass Dense layer → Down / m/p / up
γ, β.

d  Hour glass Fusion layer ← Bidirectional SSM  ⟶ ≈ Latm.

Down | Down  ⟶ Scale

m/p | m/p

m/p

up  decoder ( m/p Calculation)

□ □ . . . □ □ □

⟶ Prediction ( noise, covariance).