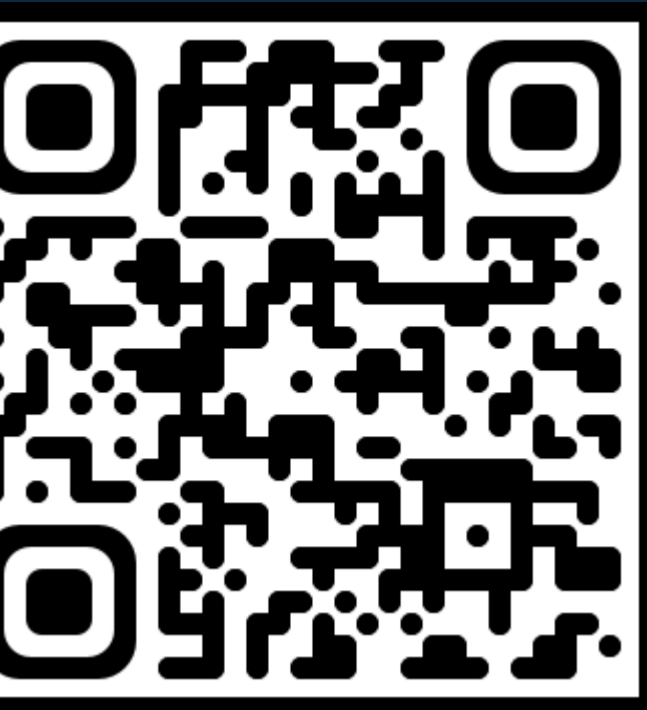


Reviving ConvNeXt for Efficient Convolutional Diffusion Models

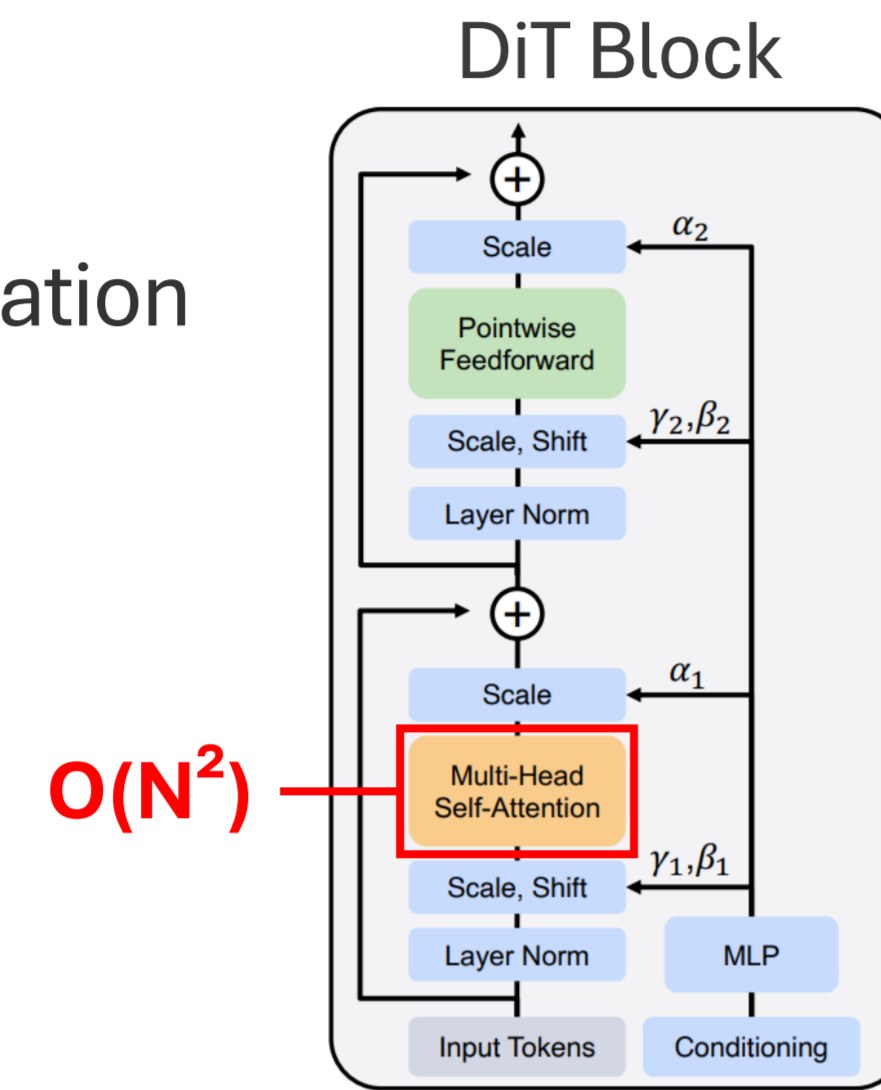
Taesung Kwon, Lorenzo Bianchi, Lennart Wittke, Felix Watine, Fabio Carrara, Jong Chul Ye, Romann Weber, Vinicius Azevedo
KAIST · ETH Zürich · ISTI-CNR · University of Pisa



Can we generate high-quality images using only convolutions? $O(N^2) \rightarrow O(N)$

Motivation

- Diffusion Transformers (DiTs) dominate image generation
- Self-attention is $O(N^2)$ — heavy compute cost
- **Goal: $O(N^2) \rightarrow O(N)$ while keeping quality**



Background: ConvNeXt

- ConvNeXts proved that modernized ConvNets can match ViT in classification
- High performance is unlocked via key designs: inverted bottlenecks, large 7x7 separable convolutions, and global response normalization
- **Can ConvNeXt's efficiency and strengths transfer to generative diffusion?**

Key Contributions

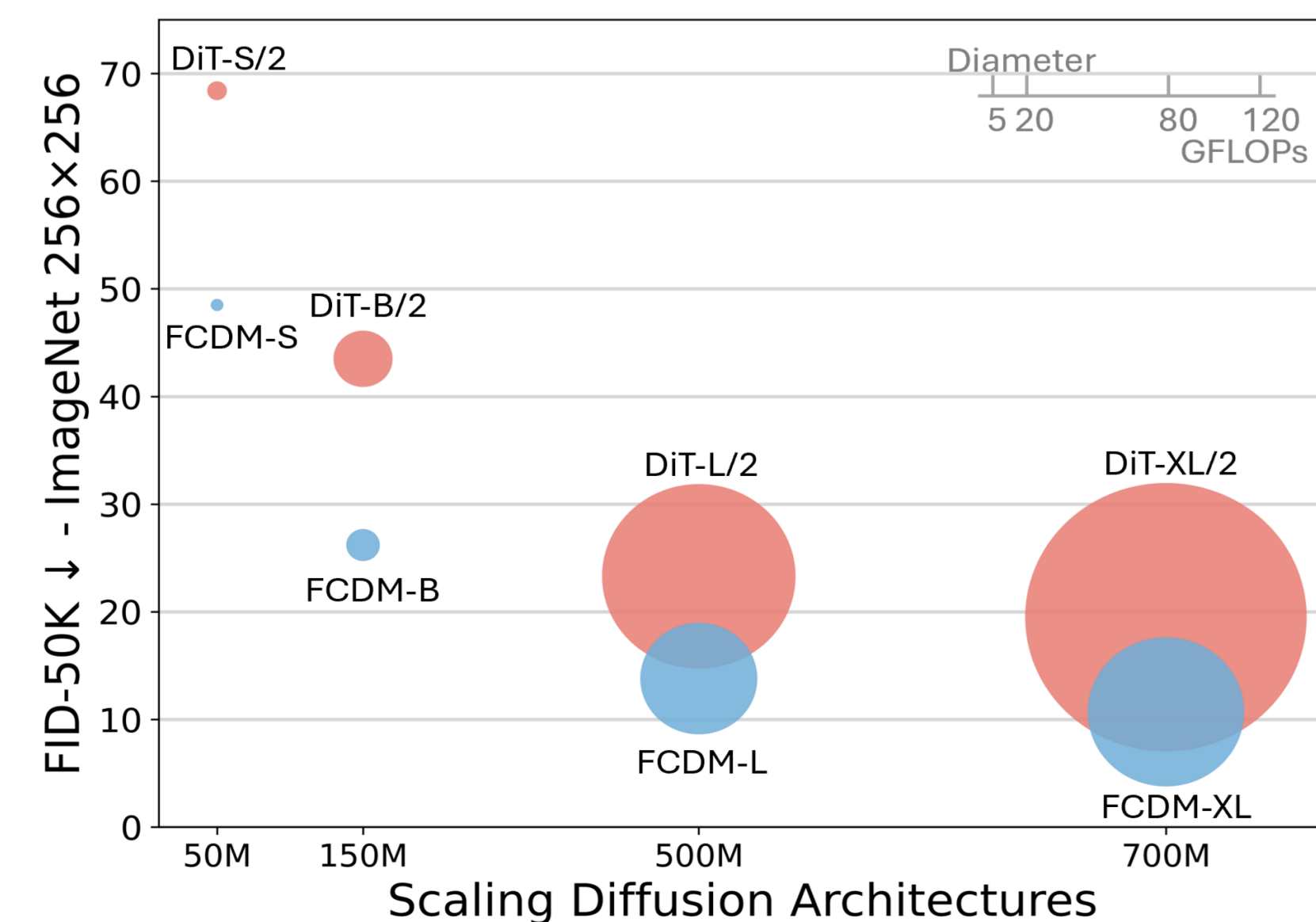
1. First fully convolutional diffusion backbone (FCDM) based on ConvNeXt architecture
2. Competitive FID with DiT-XL/2 at 50% compute cost (FLOPs)
3. 7× faster FID convergence demonstrating training efficiency

Superior Efficiency

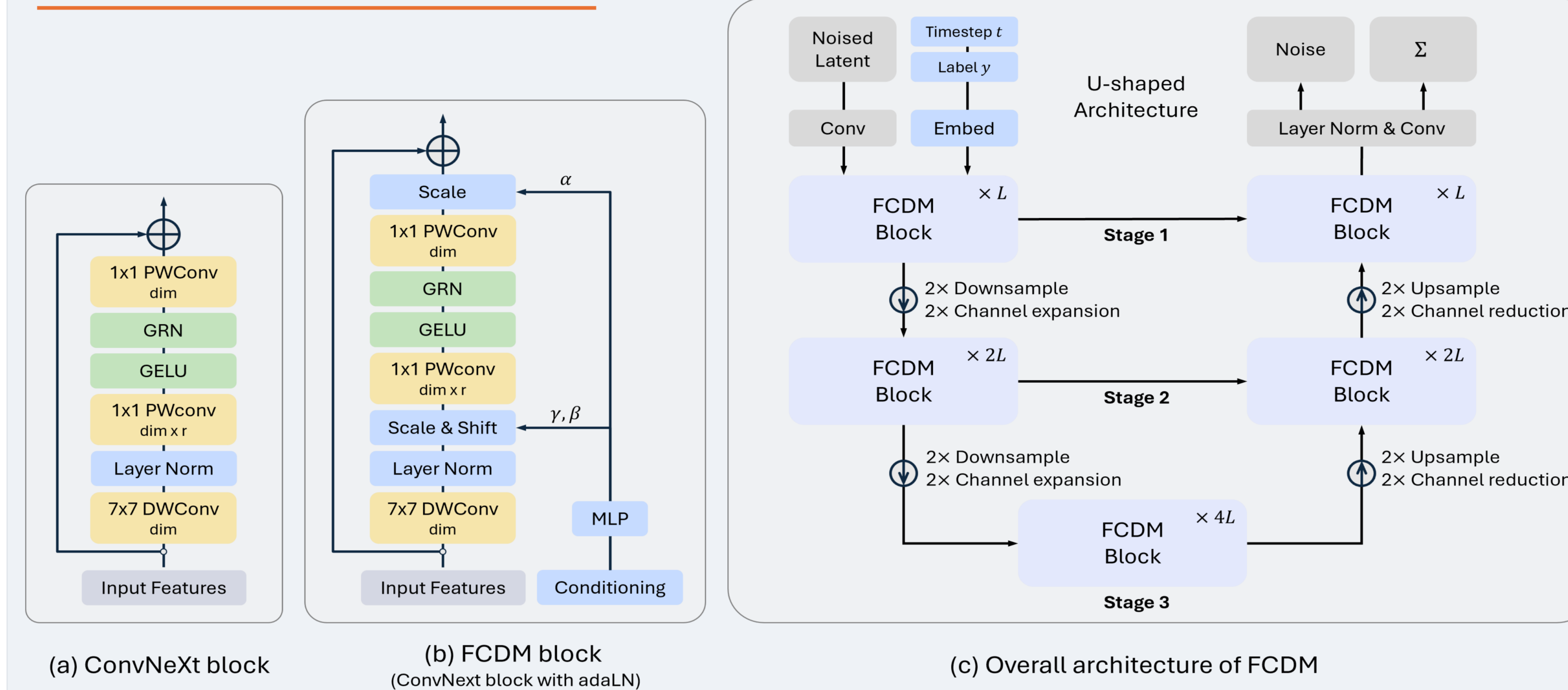
50%
Fewer FLOPs

7×
Faster Training

4 - 4090 GPUs
Only Needed

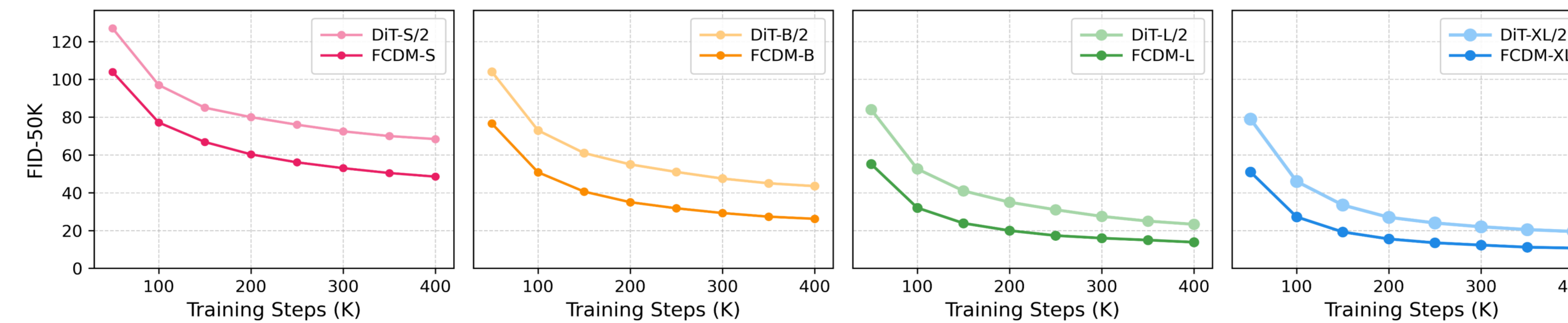


FCDM Architecture

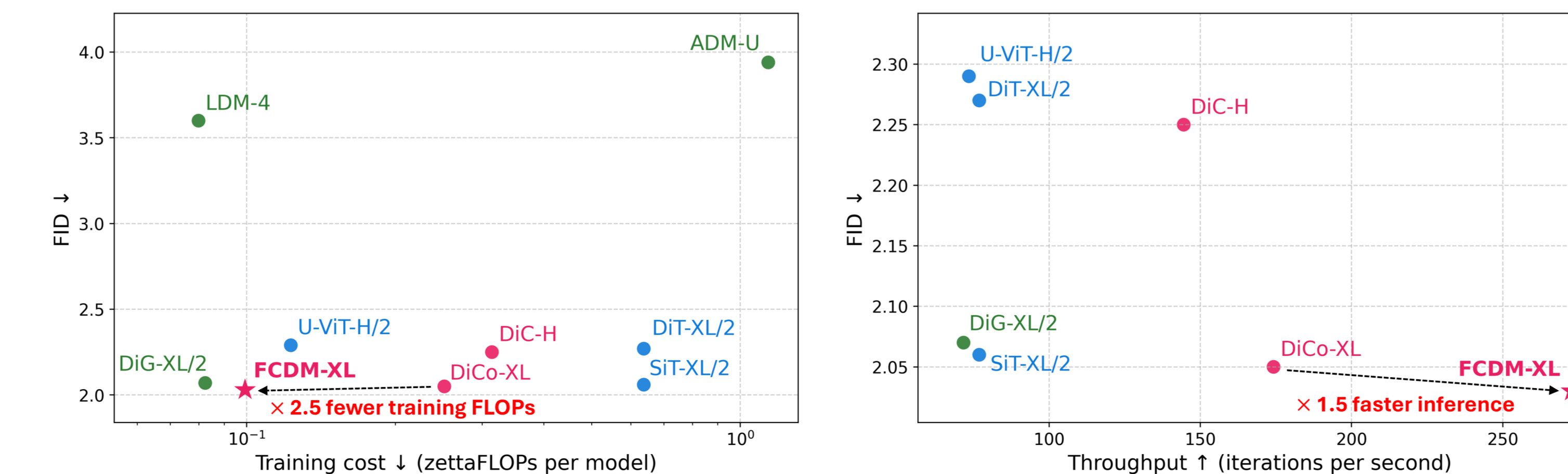


- ConvNeXt block + Adaptive LayerNorm for conditioning
- Easy scalable U-shaped design from S to XL variants
- Fully convolutional — no self-attention layers

Training Convergence



Efficiency Comparison



Generated Samples (ImageNet 512x512 & 256x256)



Text-to-image (MS-COCO 256x256)

