# "Principal Components" Enable A New Language CVPR// of Images by Xin Wen\*, Bingchen Zhao\*, Ismail Elezi, Jiankang Deng, and Xiaojuan Qi GMCV Workshop

# What makes this tokenizer different?

- 1D Causal Tokenization An ordered structure where token importance follows a hierarchical pattern.
- PCA-Like Structure Earlier tokens contain most significant information, while later tokens refine details.
- Semantic-Spectrum Decoupling ensuring tokens capture high-level rather than low-level artifacts.



**TiTok: both semantic** details and spectral

> Semantic clarity w/ spectrum across token counts.

hierarchy, mirroring

global precedence effect in human

Method	#Token	Dim.	VQ	rFID↓	<b>PSNR</b> ↑	SSIM↑	Gen. Model	Type	#Token	#Step	gFID↓	IS↑
MaskBit [55]	256	12	1	1.61	—	—	MaskBit	Mask.	256	256	1.52	328.6
RCG (cond.) [27]	1	256	X	_	_	_	MAGE-L	Mask.	1	20	3.49	215.5
MAR [28]	256	16	X	1.22	—	_	MAR-L	Mask.	256	64	1.78	296.0
TiTok-S-128 [60]	128	16	1	1.71	_	-	MaskGIT-L	Mask.	128	64	1.97	281.8
TiTok-L-32 [60]	32	8	1	2.21	_	-	MaskGIT-L	Mask.	32	8	2.77	194.0
VQGAN [13]	256	16	1	7.94	—	-	Tam. Trans.	AR	256	256	5.20	280.3
ViT-VQGAN [57]	1024	32	1	1.28		1	VIM-L	AR	1024	1024	4.17	175.1
RQ-VAE [26]	256	256	1	3.20	—	_	RQ-Trans.	AR	256	64	3.80	323.7
VAR [49]	680	32	1	0.90		-	VAR- <i>d</i> 16	VAR	680	10	3.30	274.4
ImageFolder [29]	286	32	1	0.80		_	VAR- <i>d</i> 16	VAR	286	10	2.60	295.0
LlamaGen [48]	256	8	1	2.19	20.79	0.675	LlamaGen-L	AR	256	256	3.80	248.3
CRT [39]	256	8	1	2.36	_	-	LlamaGen-L	AR	256	256	2.75	265.2
Causal MAR [28]	256	16	X	1.22	—	-	MAR-L	AR	256	256	4.07	232.4
SEMANTICIST w/ DiT-L	256	16	X	0.78	21.61	0.626	<i>ϵ</i> LlamaGen-L	AR	32	32	2.57	260.9
Semanticist w/ DiT-XL	256	16	×	0.72	21.43	0.613	€LlamaGen-L	AR	32	32	2.57	254.0

### Competitive understanding, reconstruction, and AR generation. Top: quantitative results on ImageNet reconstruction and autoregressive generation. Left: linear probing; Right: reconstruction FID.







香港大學 THE UNIVERSITY OF HONG KONG













# 5. Discussion

### **Potential Applications**

#### Limitations/Future Work

- Inference Speed: Diffusion decoding is slower than direct pixel regression.
- Alternative Architectures: Flow-matching or consistency models could improve efficiency.
- Random Ordering: Compatible to random-ordered tokens could be useful.

![](_page_0_Picture_39.jpeg)

![](_page_0_Picture_40.jpeg)

![](_page_0_Picture_41.jpeg)

# 4. AR Generation with A Variable Number of Tokens (w/LlamaGen)

![](_page_0_Picture_43.jpeg)

Models

![](_page_0_Picture_44.jpeg)

• Unified understanding & generation, image compression, data analysis, etc.