



<https://mp.weixin.qq.com/s/miEziKZNdlnRym0qJlvqLw>

- 1.
- 2. NN
- 3. FP16
- 10
- INT8 FP8 FP16 FP4
- 4. NLP SOTA Trans former 750x/2yrs
- 410x/2yrs FLOPS
- 3.0x/2yrs DRAM 1.6x/2yrs 1.4x/2yrs

DRAM

[https://people.inf.ethz.ch/omutlu/pub/Ramulator2_arxiv23.p](https://people.inf.ethz.ch/omutlu/pub/Ramulator2_arxiv23.pdf)

df

KAN MLP

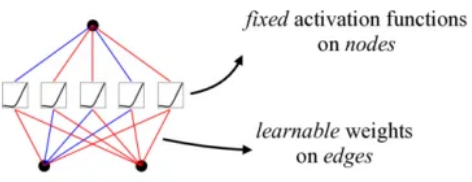
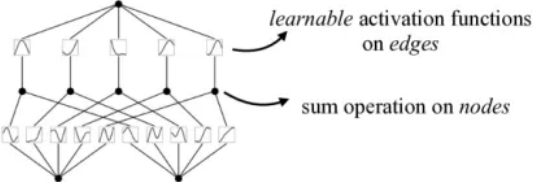
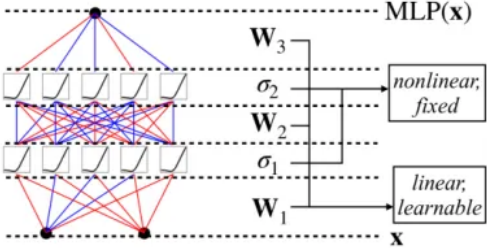
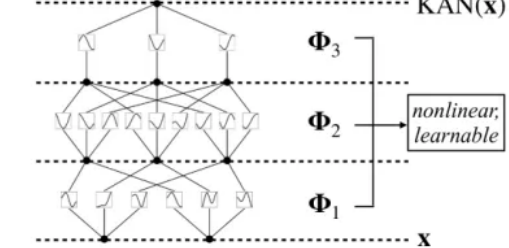
<https://mp.weixin.qq.com/s/Hrp5v5enYlx3cVwtG63d6w>

KAN MLP

Kolmogorov-Arnold

KAN

MLP

Model	Multi-Layer Perceptron (MLP)	Kolmogorov-Arnold Network (KAN)
Theorem	Universal Approximation Theorem	Kolmogorov-Arnold Representation Theorem
Formula (Shallow)	$f(\mathbf{x}) \approx \sum_{i=1}^{N(e)} a_i \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i)$	$f(\mathbf{x}) = \sum_{q=1}^{2n+1} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right)$
Model (Shallow)	<p>(a)</p> 	<p>(b)</p> 
Formula (Deep)	$\text{MLP}(\mathbf{x}) = (\mathbf{W}_3 \circ \sigma_2 \circ \mathbf{W}_2 \circ \sigma_1 \circ \mathbf{W}_1)(\mathbf{x})$	$\text{KAN}(\mathbf{x}) = (\Phi_3 \circ \Phi_2 \circ \Phi_1)(\mathbf{x})$
Model (Deep)	<p>(c)</p> 	<p>(d)</p> 

Revision #1

Created 11 January 2025 09:44:04 by Colin

Updated 11 January 2025 09:44:05 by Colin