

?? ??? Binary ????????

??

1. FPGA [] LUT([]) [] LUT []
 1. []
 2. []
2. FPGA [] AI [] PPA
 1. []
 2. [] FPGA []
 3. [] PPA []
3. [] FPGA [] LUT [] AI []
 1. []
4. []
 1. [] []
 2. [] LUT []
5. binary LUT [] bit []
 1. [] 8bit [] 256bit []
 2. []
 3. [] CIM []
6. [] Bitnet []
 1. [] FP8 [] int4 []
 2. []

??

1. LUT6
 1. 6 [] 1 [] LUT [] 64 [] 64bit
 2. xc7k480t [] 74659 slices [] slice [] 4 LUT [] 8 [] flip-flop
 3. [] **lut** [] [] 64bit*74659*4 = 2389088 byte [] 74659*4
 4. [] LUT [] LUT []
2. []
 1. []
 2. 4 [] 6LUT [] 8bit lut [] 8 [] LUT8 [] 8bit [] 8 [] 8 []
 3. []
3. [] LUT []
 1. [] LUT []

1. "[input] + weight" [output] LUT
2. [input] LUT [output] LUT
 - [input]
 - LUT [output] LUT [output]
1. [input] GPU [output]
2. [input]
3. [input] AI [output]
 - FP8 [output]
2. [input] [output] [output]
4. [input] LUT [output] $f(x)=y$ [output] LUT [output]
 1. ElementWise [output] CNN [output] GEMM [output] weight [output] LUT [output]
 1. [input] LUT [output]
 2. Reduction [output] normal [output] softmax [output] Reduce [output]
 3. [input] GEMM [output]
5. [input] LUT
 1. LUT [output] FPGA [output]

?????

????

1. [input] "[input]" [output]
2. [input] softmax [output] normal [output] GEMM [output] CNN [output]
3. [input] LUT [output]
4. [input]

??????

1. [input] LUT [output]
 1. [input]
 2. [input] "[input]" [output]
 - [input]
 3. [input]
 - [input]
2. [input]
 1. [input]
 2. [input]
 1. [input]
 3. [input]
 1. [input] bit [output] XOR [output] AND [output]
 4. [input] LUT [output]
 - [input]
3. [input]
 1. [input] linear [output]

1. $\text{input} \times \text{weight} = \text{output}$
2. $\text{input} + \text{weight} = \text{output}$ LUT
2. $\text{input} \Rightarrow \text{output}$

AI???????

1. input
2. weight
 1. $\text{input} \times \text{weight}$
3. output
4. $\text{input} + \text{weight}$
5. input
6. weight
7. input shift scale

LLM??

1. $\text{input} \times \text{weight} = \text{output}$ bit
 1. input
2. weight
 1. $\text{input} \times \text{weight}$
 2. $\text{input} + \text{weight}$
 3. $\text{input} \Rightarrow \text{output}$
3. input 8bit
 1. input bit
 2. 8bit \times 256 weight
 3. weight bit
 1. $\text{input} \times \text{weight}$
 4. output
 1. input LUT
 2. input LUT bit LUT
 3. input
4. input
 1. $\text{input} / \text{weight}$
 2. $\text{input} \times \text{weight}$ " " " " "
 3. input
 1. input
 2. input
 3. input
5. input
 1. $x \rightarrow x$
 1. input attention qkv bit
 1. QK input qk input input
 2. input LUT reduce
6. input LUT
 1. input

9. Xilinx UltraScale Architecture Configurable Logic Block User Guide
https://japan.xilinx.com/support/documentation/user_guides/ug574-ultrascale-clb.pdf
 10. 1-bit AI Infra: Part 1.1, Fast and Lossless BitNet b1.58 Inference on CPUs
<https://arxiv.org/abs/2410.16144v2>
 11. Bitnet.cpp: Efficient Edge Inference for Ternary LLMs <https://arxiv.org/abs/2502.11880v1>
 12. Continual Quantization-Aware Pre-Training: When to transition from 16-bit to 1.58-bit pre-training for BitNet language models? <https://arxiv.org/abs/2502.11895v1>
 13. (NEW!) BitNet v2: Native 4-bit Activations with Hadamard Transformation for 1-bit LLMs
<https://arxiv.org/abs/2504.18415>
 14. (NEW!) BitVLA: 1-bit Vision-Language-Action Models for Robotics Manipulation
<https://arxiv.org/abs/2506.07530>
-

Revision #71

Created 2025-04-12 08:03:12 UTC by Colin

Updated 2026-04-29 07:33:44 UTC by Colin