

Binary AI

- [\[\] \[\] Binary \[\] \[\] \[\] \[\] \[\] \[\]](#)
- [\[\] \[\] \[\] \[\]](#)

?? ??? 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 / 6 * 2 * 0.5G GFlops = 50 TFlops
 4. LUT LUT
2.
 - 1.
 2. 4 6LUT 8bit lut 8 LUT8 8bit 8 8
 - 3.
3. LUT
 1. LUT
 1. " +weight" LUT

2. LUT LUT

LUT LUT
1. GPU
2.
3. AI
FP8

2.
4. LUT $f(x)=y$ LUT
1. ElementWise CNN GEMM weight LUT
1. LUT
2. Reduction normal softmax Reduce
3. GEMM
5. LUT
1. LUT FPGA

?????

????

1. " "
2. softmax normal GEMM CNN
3. LUT
4.

??????

1. LUT
1.
2. " "
3.
2.
1.
2.
1.
3.
1. bit XOR AND
4. LUT
3.
1. linear
1. = =

- 2. int “ int +weight” int LUT
- 2. int \Rightarrow int

AI???????

- 1. int
- 2. int
 - 1. int
- 3. int
- 4. int
- 5. int
- 6. int
- 7. int shift int scale int

LLM??

- 1. int “ int ” “ int bit”
 - 1. int
- 2. int
 - 1. int
 - 2. int “ int ” “
 - 3. int int int
- 3. int 8bit int
 - 1. int bit int
 - 2. 8bit int 256 int weight
 - 3. weight int bit int
 - 1. int
- 4. int
 - 1. int LUT int
 - 2. int LUT int bit int LUT
 - 3. int
- 4. int
 - 1. int / int
 - 2. int “ int ” “ int ” “ int ” “ int ”
 - 3. int
 - 1. int
 - 2. int
 - 3. int
- 5. int
 - 1. $x \rightarrow x$
 - 1. int attention int qkv int bit int
 - 1. QK int qk int input int
 - 2. int LUT int reduce
- 6. int LUT int
 - 1. int

??

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>

???

????LLM?Binary

- 1. Lut[] touch[]
- 2. [] Lut[]

?????Bianry??

- 1. []
- 1. []