

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

??

- 1. FPGA LUT() LUT
 - 1.
 - 2.
- 2. FPGA AI
 - 1. PPA
 - 2. FPGA
 - 3. PPA
- 3. FPGA LUT AI
 - 1.
- 4.
 - 1.
 - 2. LUT

??

- 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
 - 3.
- 3. LUT
 - 1. LUT
 - 1. “+weight” LUT
 - 2. 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. ☐ " ☐ " ☐
2. ☐
 1. ☐
 2. ☐
 1. ☐
 3. ☐
 1. ☐ bit ☐ XOR ☐ AND ☐
 4. ☐ LUT
3. ☐
 1. ☐ linear = ☐ = ☐
 1. ☐ = ☐
 2. ☐ " ☐ +weight" ☐ LUT
 2. ☐ => ☐

AI??????

1. ☐
2. ☐
 1. ☐
3. ☐
4. ☐
5. ☐

6. ☐
7. ☐ shift ☐ scale ☐

LLM??

1. ☐ " ☐ " ☐ bit ☐
 1. ☐
2. ☐
 1. ☐
 2. ☐ " ☐ "
 3. ☐ ☐ ☐
3. ☐ 8bit ☐
 1. ☐ bit ☐
 2. 8bit ☐ 256 ☐ weight
 3. weight ☐ bit ☐
 1. ☐
4. ☐
 1. ☐ LUT ☐
 2. ☐ LUT ☐ bit ☐ LUT
 3. ☐
4. ☐
 1. ☐ / ☐
 2. ☐ " ☐ " ☐ " ☐ "
 3. ☐
 1. ☐
 2. ☐
 3. ☐
5. ☐
 1. $x \rightarrow x$
 1. ☐ attention ☐ qkv ☐ bit ☐
 1. QK ☐ qk ☐ input ☐
 2. ☐ LUT ☐ reduce
6. ☐ **LUT** ☐
 1. ☐

??

1. ☐ LLM/CNN
 1. ☐
 2. ☐ PPA
 3. ☐ FPGA ☐
2. ☐

???

1. ☐
2. AI ☐

3.
4. DNA
5.
6. LUT
1.
2.
3. scale

??????

1. BitNet b1.58 <https://mp.weixin.qq.com/s/G9ZbMnBVbeH1m45HY2JlKA>
2.

??

1.
2.

??

1. <https://arxiv.org/html/2502.19008v1>
2. Binary Neural Networks <https://zhuanlan.zhihu.com/p/117285043>
3. <https://github.com/ryuz/BinaryBrain>
4.
<https://www.slideshare.net/kentotajiri/ss-77136469>
5. BinaryConnect: Training Deep Neural Networks with binary weights during propagations
<https://arxiv.org/pdf/1511.00363.pdf>
6. Binarized Neural Networks
<https://arxiv.org/abs/1602.02505>
7. Binarized Neural Networks: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1
<https://arxiv.org/abs/1602.02830>
8. XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks
<https://arxiv.org/abs/1603.05279>
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 #66

Created 12 April 2025 08:03:12 by Colin

Updated 7 July 2025 16:56:24 by Colin