




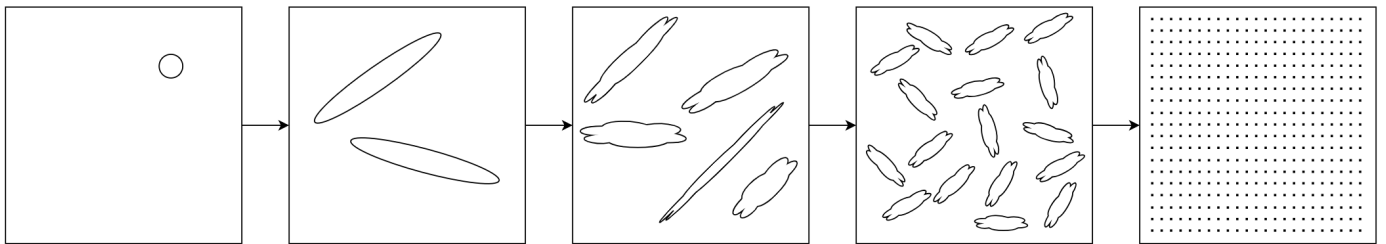











LLM????????

??

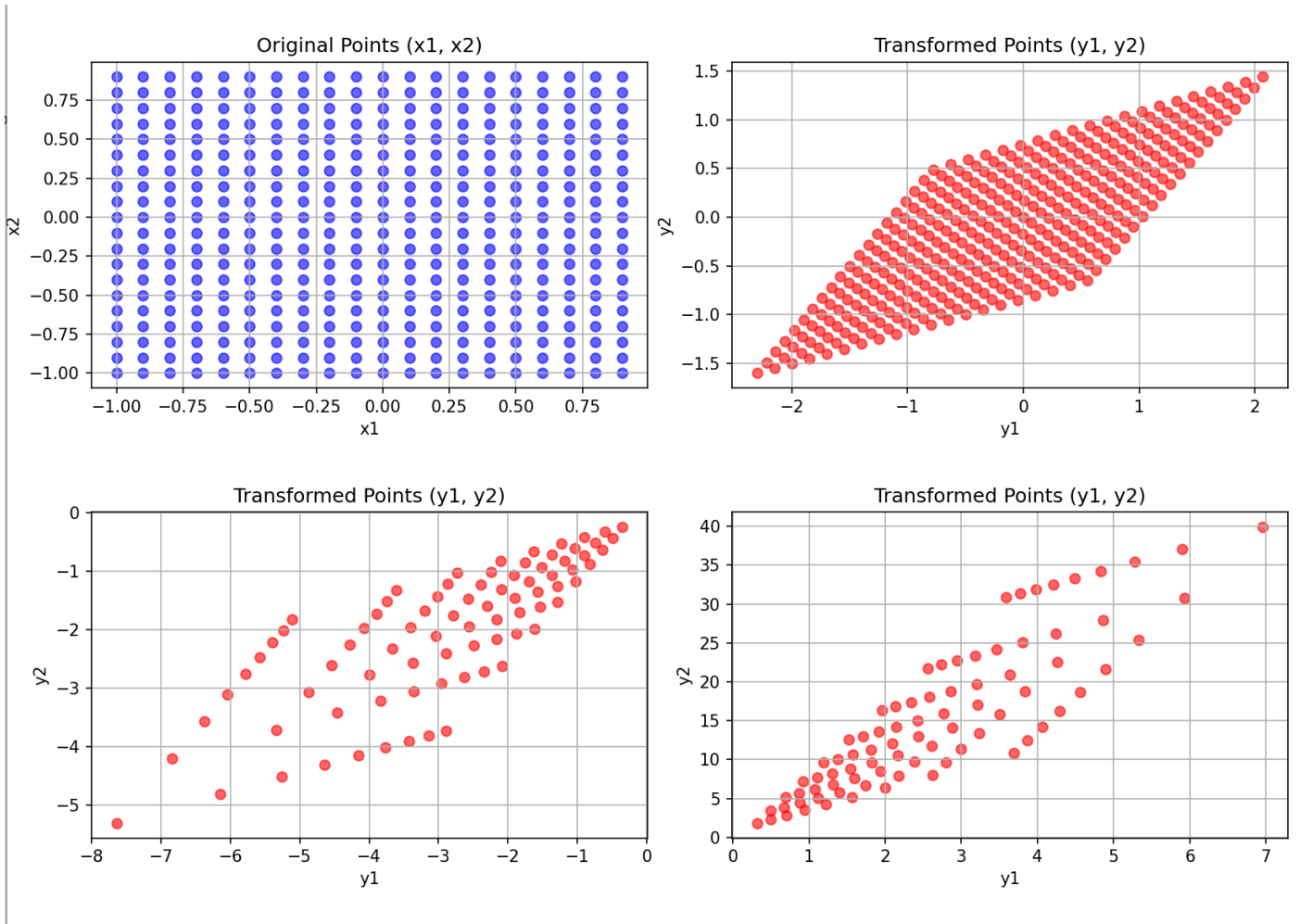
1. 32 bit  32 bit 
2. bit 
3.  32 

??????????



1.  hidden status 
2. 
 1. 
 2. **hidden status** 
3.  index
4.  DeekseekV3 671B  98%  256×60
 $*44\text{MB}=670\text{GB}$
5. $\text{expert} \times \text{linear} = \text{hidden_size} * \text{moe_intermediate_size} * 3 = 7168 * 2048 * 3 = 44\text{MB}$

???? ??????



???Dot

- $$A = B * (C + D) \quad A = B * C + B * D$$
 -
- ResNet
 - $$A = B * C + B$$

LLM?????

-
- LLM
 -
 -
- Plinko
 -
 - layer token hidden status
 - hidden status token
- token
 - LLM +
 -

Revision #18

Created 2025-07-27 06:46:14 UTC by Colin

Updated 2026-04-29 07:34:27 UTC by Colin