**Decode Time per Token (4 GPU Average)**

Legend:
- Deepseek-16b
- Deepseek-v2-lite
- Qwen1.5-moe
- Qwen3-30B

X-axis: Decode Token Index
Y-axis: Time per Token (s)