**(a) Cross-Platform Inference Latency**

- A100 (Cloud): 2.06ms (250W)
- Jetson Orin (Edge): 10.69ms (15W)
- Ascend 310 (Edge): 6.37ms (8W)
- Hygon DCU (Domestic): 8.22ms (150W)

**(b) A100 Throughput vs Batch Size**

| Batch Size | Throughput (FPS) |
|---|---|
| $2^0$ | 486 |
| $2^2$ | 1911 |
| $2^3$ | 3770 |
| $2^4$ | 7540 |
| $2^5$ | 8748 |
| $2^6$ | 9678 |