

F1 Score

$2 \times 10^1$

$10^1$

- Quantization
- Full Cache
- InfiniGen
- KVDrive

Memory Size (GB)

1

2

3

4

