| Type | Optimization | Threat | FT 22 | Orca 22 | vLLM 23 | FlexGen 23 | FastServe 23 | FastServe2 23 | Sarathi 23 | Lookahead 24 | REST 24 | SpecInfer 24 | Medusa 24 | DistServe 24 | Splitwise 24 | LoongServe 24 | TetriInfer 24 | InfiniteLLM 24 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Batch** | Iteration-Level Batch | | | Initial | ✓ | | ✓ | | ✓ | | | | | ✓ | | | ✓ | |
| | Chunked Prefill | | | | | | | | Initial | | | | | | | | ✓ | |
| | Prepack Prefill | | | | | | | | | | | ✓ | | ✓ | | | ✓ | |
| **Parallel** | Speculation | S | | | | | | | | ✓ | ✓ | ✓ | ✓ | | | | | |
| | Prompt-Based Speculation | S | | | | | | | | ✓ | | | | | | | | |
| | Context-Based Speculation | S | | | | | | | | | ✓ | | | | | | | |
| | Tensor Parallelism | | ✓ | | | | | | | | | | | | | | | |
| | SafeTensors | | | | | | | | | | | | | | | | | |
| | Sequence Parallelism | | | | | | | | | | | | | | | ✓ | | |
| **Memory** | Paging | | | | Initial | | | | ✓ | | | | | | | | ✓ | |
| | Multi-Query Attention | | | | | | | | | | | | | | | | | |
| | Grouped-Query Attention | T | | | | | | | | | | | | | | | | |
| **Tranmission** | Offloading | SE | | | ✓ | ✓ | | | | | | | | | | | | |
| | Duplication | T | | | | | | | | | | | | | | | | |
| | Pulling | SET | | | | | | | | | | | | ✓ | | | | |
| | Request Migration | | | | | | | | | | | | | | | ✓ | | |
| | Disaggregated Arch | | | | | | | | | | | | | ✓ | ✓ | | ✓ | |
| **Scheduling** | Priority-Based | T | | | | | ✓ | | ✓ | | | | | | | | ✓ | |
| | Request-Level Prediction | T | | | | | ✓ | | | | | Small Model | | | | | ✓ | |
| | Machine-level Scheduler | ET | | | | | ✓ | | | | | | | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Instance Flip | | | | | | | | | | | | | | ✓ | | ✓ | |
| | Global Profiling | P | | | ✓ | | | | | | | | | ✓ | ✓ | | | |
| **Verification** | Open Source | V | | | | | | | | | | | | | | | ✓ | |