Ingesting data for use with a large language model for AI:

# Latest-generation Dell™ PowerEdge™ servers powered by 5th Generation AMD EPYC™ processors offer a range of strong options
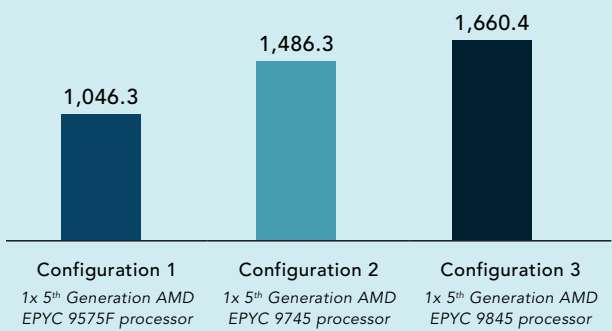
We measured the performance of multiple disaggregated infrastructure server configurations to help decision-makers choose the right one for their needs.
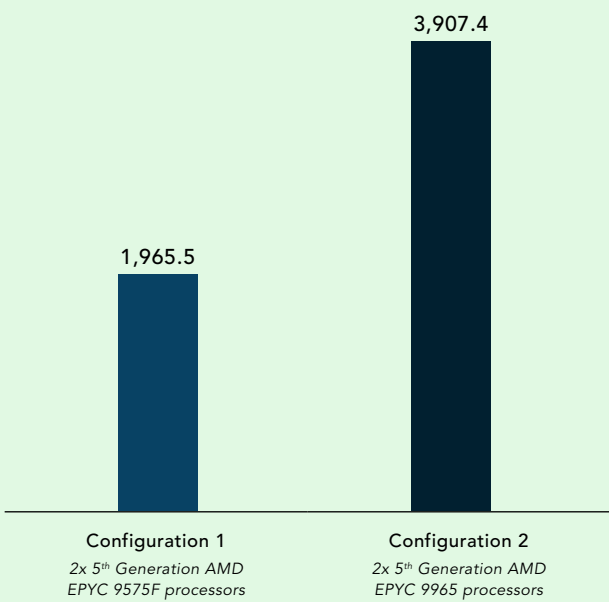
## bfloat16 precision

Up to **1,660** sentences per second
on a Dell PowerEdge R7715

Up to **3,907** sentences per second
on a Dell PowerEdge R7725

### Dell PowerEdge R7715 with bfloat16 precision
Sentences per second



| | | |
|---|---|---|
| 1,046.3 | 1,486.3 | 1,660.4 |
| Configuration 1 | Configuration 2 | Configuration 3 |
| 1x 5th Generation AMD EPYC 9575F processor | 1x 5th Generation AMD EPYC 9745 processor | 1x 5th Generation AMD EPYC 9845 processor |

### Dell PowerEdge R7725 with bfloat16 precision
Sentences per second



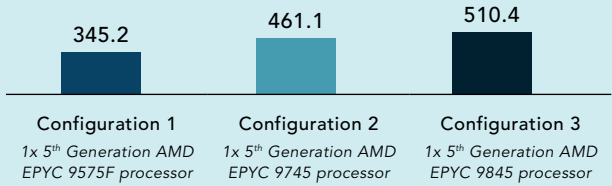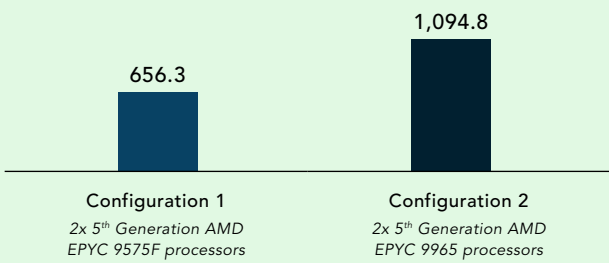| | |
|---|---|
| 1,965.5 | 3,907.4 |
| Configuration 1 | Configuration 2 |
| 2x 5th Generation AMD EPYC 9575F processors | 2x 5th Generation AMD EPYC 9965 processors |

## float32 precision

Up to **510** sentences per second
on a Dell PowerEdge R7715

Up to **1,094** sentences per second
on a Dell PowerEdge R7725

### Dell PowerEdge R7715 with float32 precision
Sentences per second



| | | |
|---|---|---|
| 345.2 | 461.1 | 510.4 |
| Configuration 1 | Configuration 2 | Configuration 3 |
| 1x 5th Generation AMD EPYC 9575F processor | 1x 5th Generation AMD EPYC 9745 processor | 1x 5th Generation AMD EPYC 9845 processor |

### Dell PowerEdge R7725 with float32 precision
Sentences per second



| | |
|---|---|
| 656.3 | 1,094.8 |
| Configuration 1 | Configuration 2 |
| 2x 5th Generation AMD EPYC 9575F processors | 2x 5th Generation AMD EPYC 9965 processors |

Configurations leveraging bfloat16 precision significantly boosted sentence processing rates, highlighting their suitability for demanding AI applications. A disaggregated architecture using these servers can help you independently scale compute and storage resources, optimizing efficiency and cost-effectiveness.

By carefully selecting the appropriate server model and processor configuration for your workload, you can achieve a balanced solution that accelerates AI ingestion while avoiding unnecessary overprovisioning, enabling faster deployment and expansion of internal AI platforms.

**Read the report** ▶

**Principled Technologies®**