

Performance Targets

- 1,000+ docs/hour
- Auto-scale 1-50 pods

Celery Architecture for Document AI Processing

Distributed Task Queue with Auto-scaling Kubernetes Workers

Auto-scaling Logic

- Queue length > 100
- CPU usage > 70%



