# Remote engines on IBM DataStage

### What's remote about them?

Here's the key idea: DataStage separates where you design pipelines (the control plane) from where they run (the data plane). The control plane lives in IBM's cloud (in Dallas, Frankfurt, or Sydney), while the data plane (the remote runtime engine) lives wherever you want.

This means:
- You get a nice web UI from IBM for building pipelines
- Your data never has to leave your infrastructure when jobs run

### Why Use a Remote Engine?

1. **Data stays put**: Process sensitive data on-premises without shipping it to IBM Cloud
2. **Avoid egress fees**: If your data lives in AWS, run the engine there instead of moving terabytes to IBM Cloud
3. **Compliance**: Meet data residency requirements by keeping processing local
4. **Extra functionality**: Use stages like Java Integration or custom functions that aren't supported in IBM Cloud

### How does it work?

The remote engine is just a container! You can run it on:
- Kubernetes (we provide manifests)
- Docker (we provide scripts)
- Really any container management platform

You deploy it once, then send DataStage jobs to it. Multiple projects can share one engine, but dedicated engines are recommended for production.

Fun fact: The engine uses outbound HTTPS connections to IBM's control plane—no inbound firewall rules needed.

### Scaling and Management

- **Elastic scaling**: Configure horizontal pod autoscaling in Kubernetes
- **Disaster recovery**: Deploy backup engines in another region
- **Observability**: Hook up Databand to monitor pipelines, in addition to IBM Cloud Monitoring

You manage the engine lifecycle, while IBM manages the control plane. IBM provides regular security patches, health monitoring dashboards, and log forwarding capabilities.

**Considerations**
- **Commitment**: Once a project uses a remote engine, you can't switch back to IBM-hosted for that project
- **Limitations**: Concurrent job capacity varies by engine size
- **Billing**: You pay for vCPU-hours of engines deployed each month
- **Network**: Requires consistent minimum 10Mbps+ connection to IBM Cloud

**Who shouldn't use this?**

If you:

- Like having all your data in IBM Cloud
- Only process small amounts of data
- Want to avoid container maintenance

Then you should probably stick with regular DataStage-aaS.

**Who should use this?**

If you:

- Have data that can't leave your infrastructure
- Want to avoid cloud transfer costs
- Want access to more connectors and stages
- Are already using CPDaaS

Then cloud-style development with on-premises data control makes remote engines worth exploring.