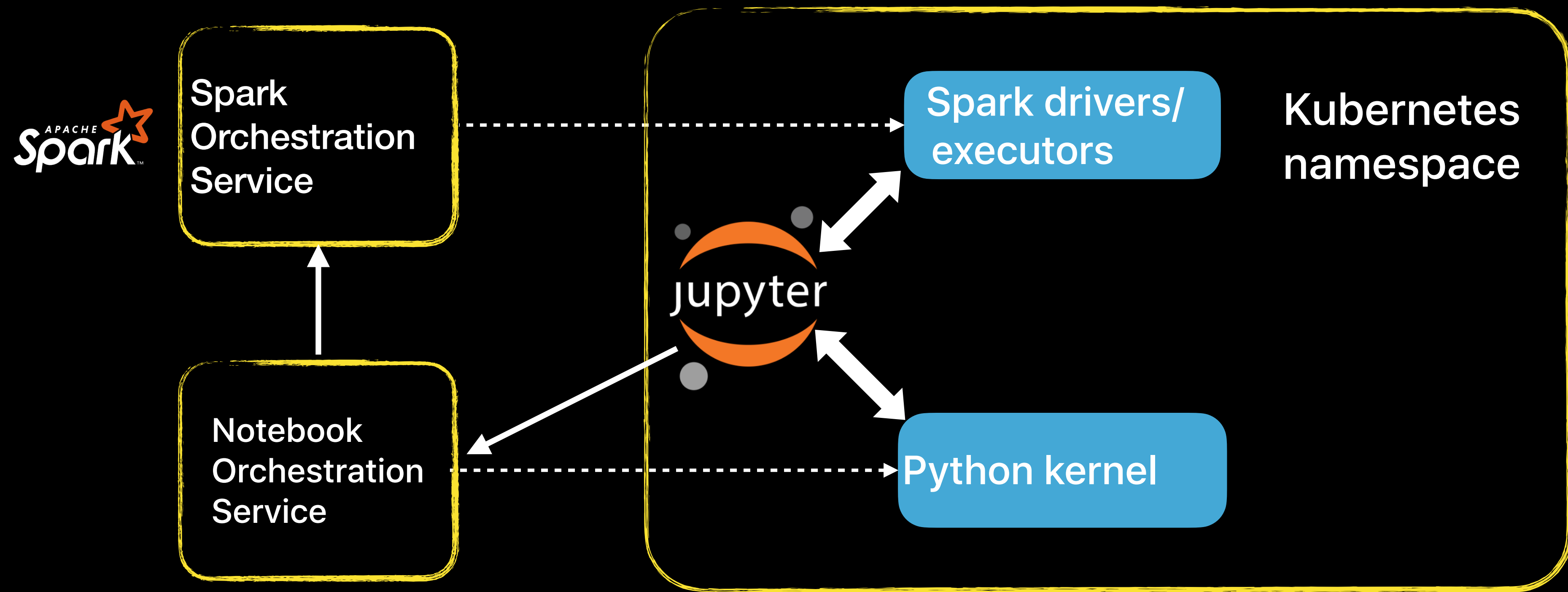# Data Science Environments In Seconds: Scaling Jupyter Notebooks in Kubernetes

Jialin Zhang

AppDeveloperCon  |  Apple  |  Nov 12, 2024
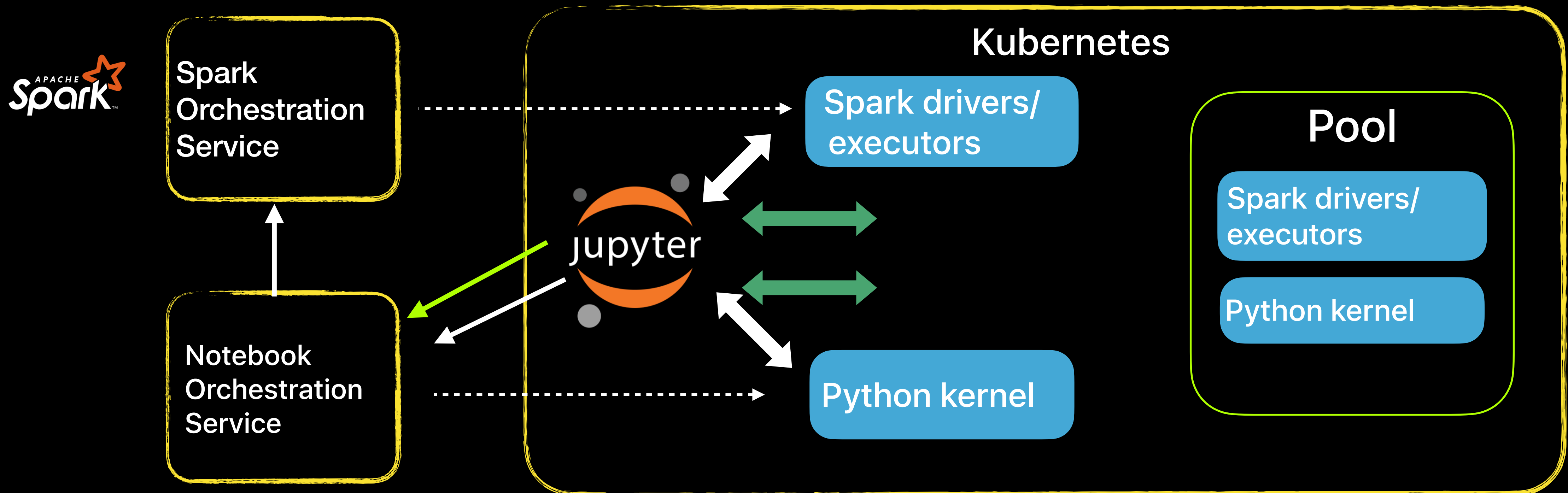
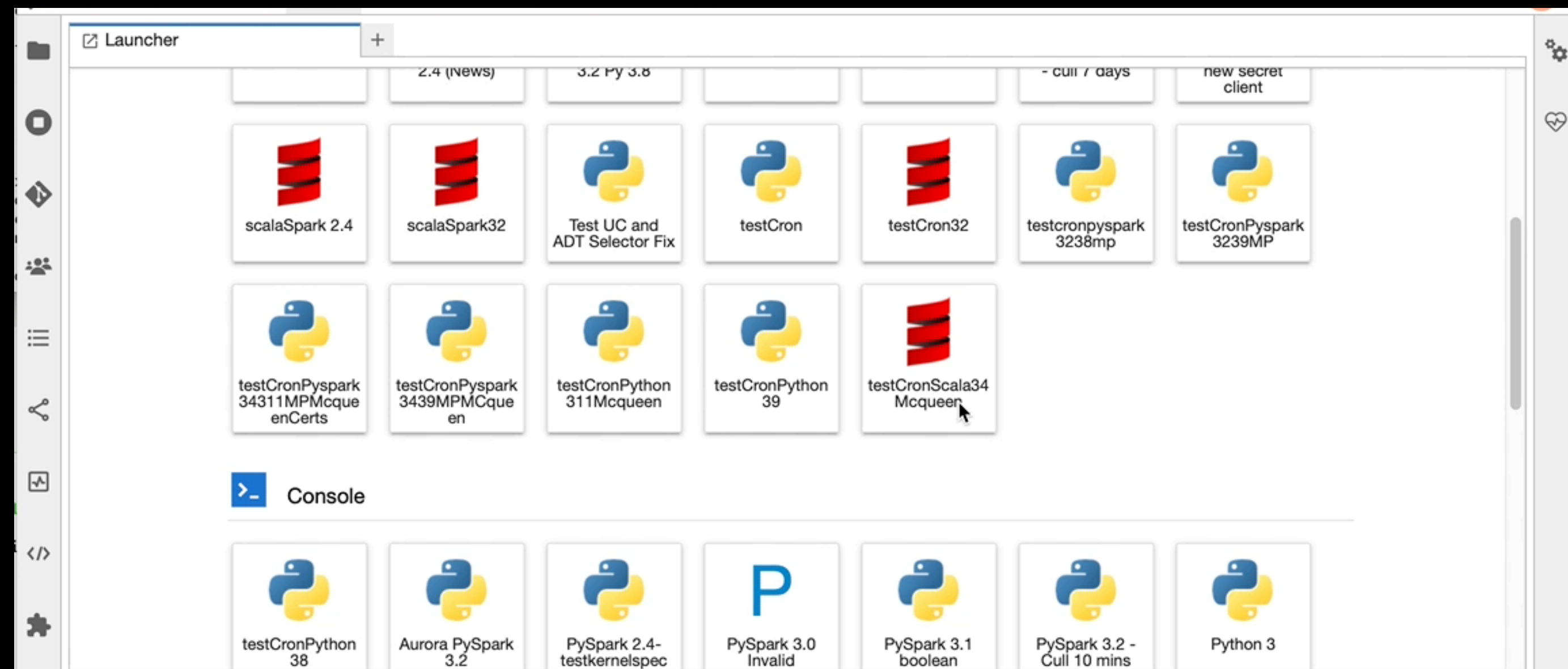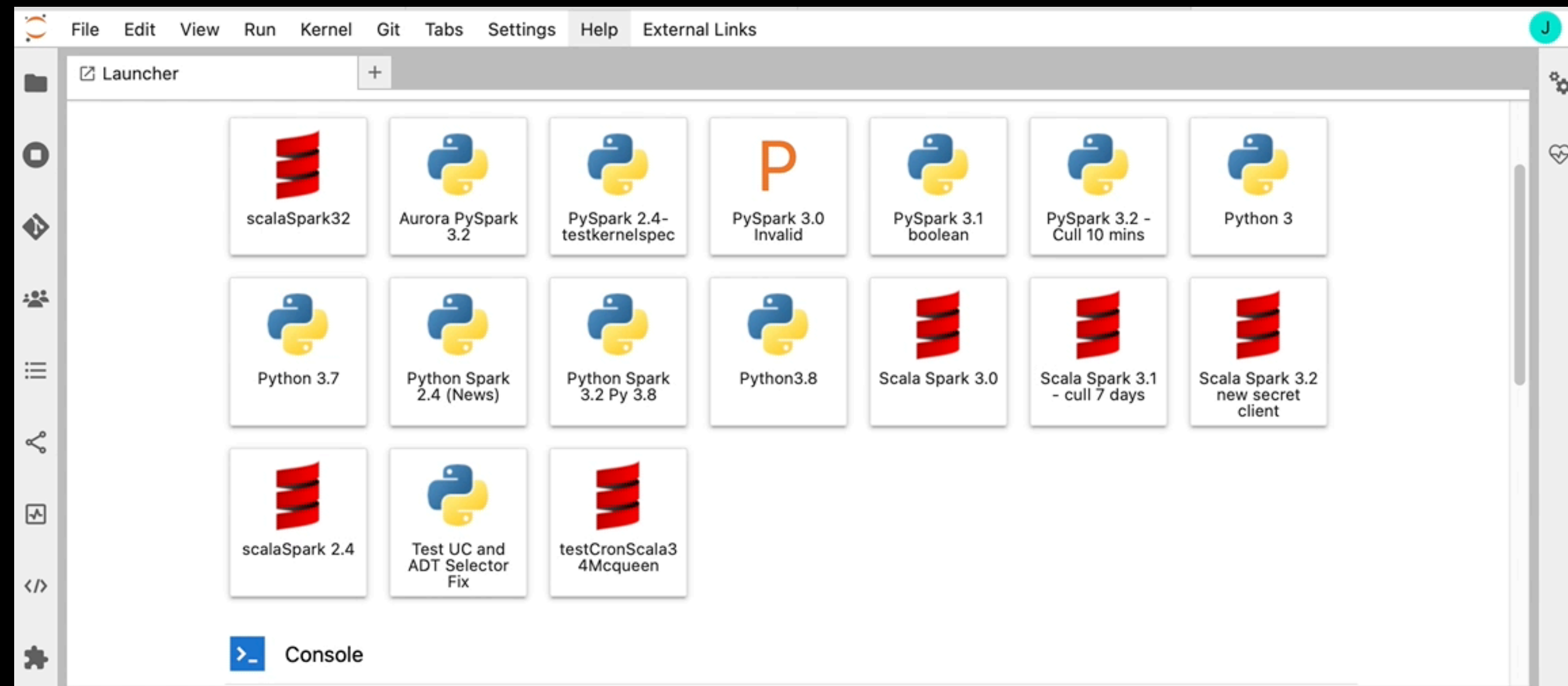# Challenge

Kubernetes Pod Kernel Launch takes minutes
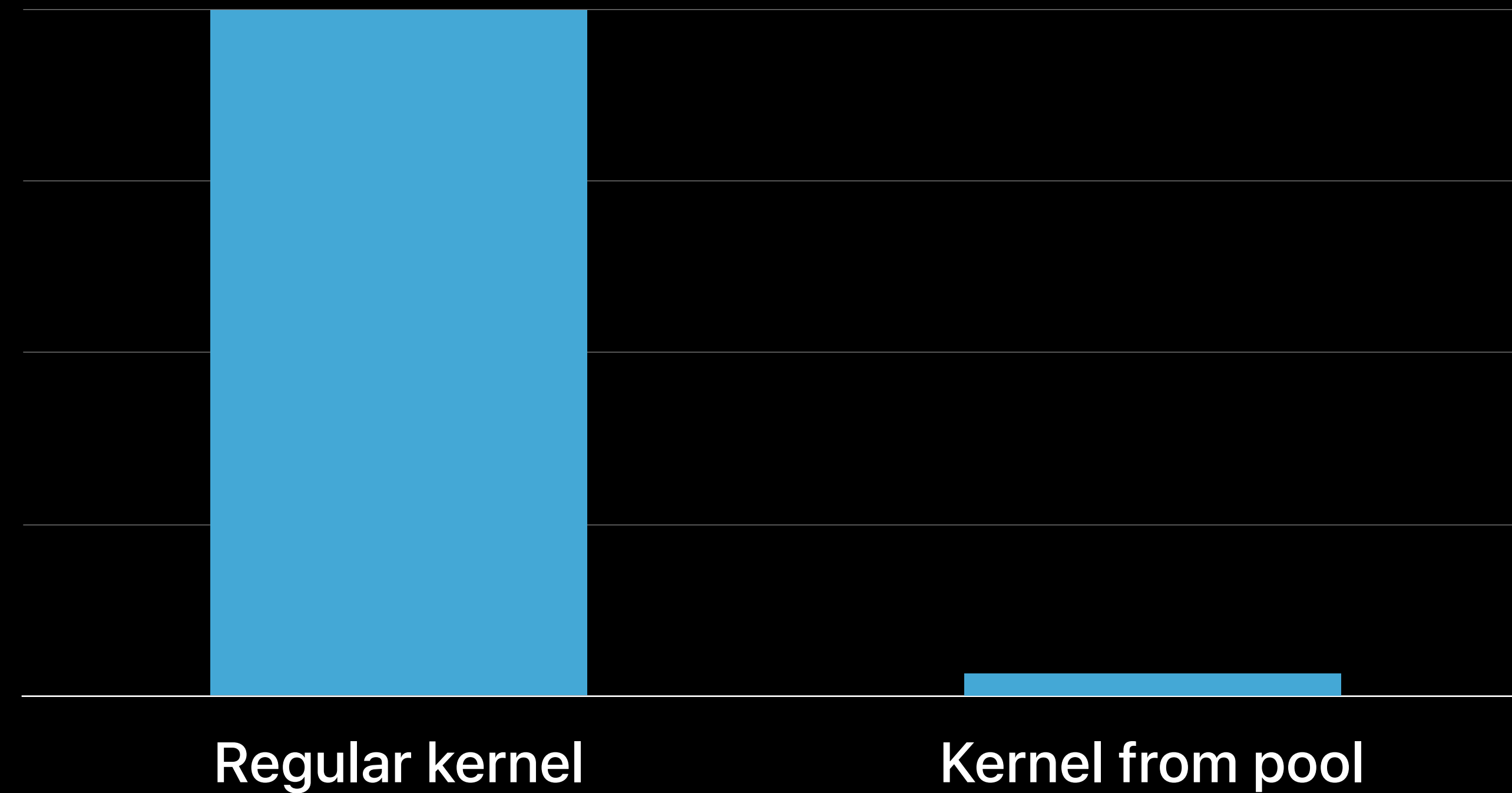
# Solution

Pre-warm kernel in pool

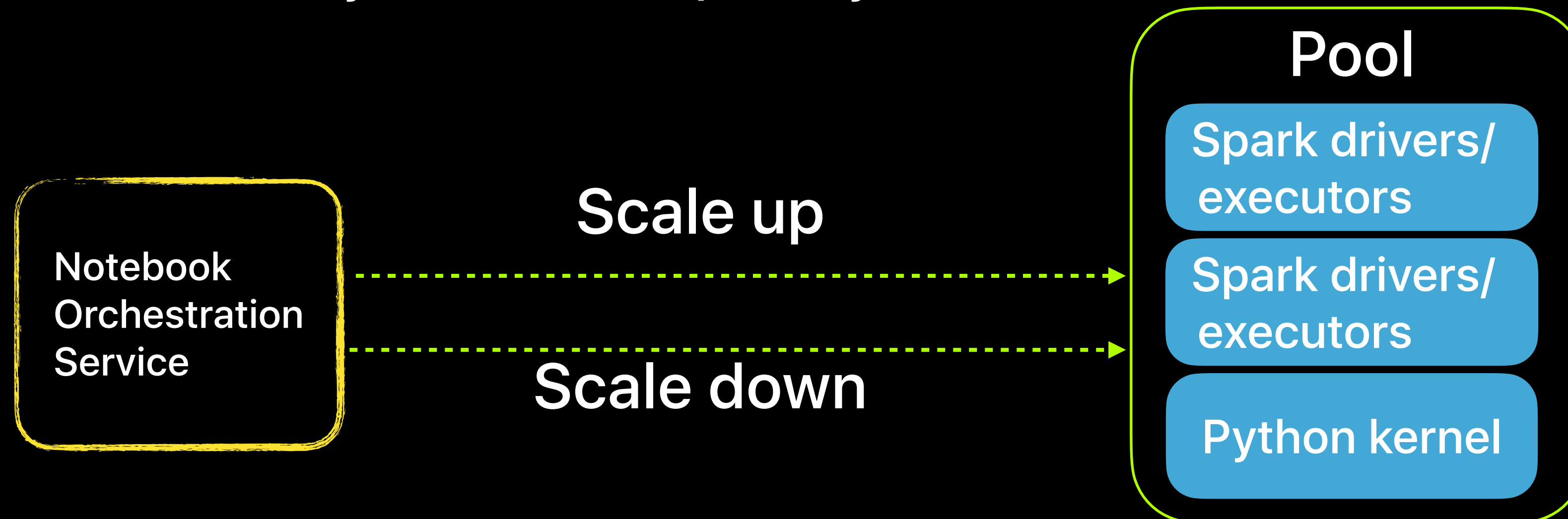If user wants kernel available in pool, directly pick up from pool

# Improved Productivity

Reduce launch latency by 96%, from minutes to seconds



Regular kernel          Kernel from pool

# Architecture + Design

Resource driven dynamic pool management

- Periodically
- Scale up/down kernels: Pool + existing usage < 95% quota usage
- Pool max size <= 10% of namespace quota
- Prediction -> rank by launch frequency



Notebook Orchestration Service

Scale up

Scale down

Pool

Spark drivers/ executors

Spark drivers/ executors

Python kernel

# Architecture + Design

## Cost Savings

- Large Namespace
- Actively detect kernels in pools that were running for a long time and tear them down
- Kernel configuration with small initial resource footprint

## Pool per Namespace

- Auto-detect eligible namespaces dynamically

# Areas to Explore

Improving prediction

Time triggered scaling —> Watcher triggered scaling

More namespaces with lower resource quota