# Productionizing

# AI + LLM Apps with Ray Serve

**Adam Breindel**
**Anyscale**

adamb@anyscale.com

@adbreind

# 👨‍💻 Meet the tutorial team!

**Marwan**

marwan@anyscale.com

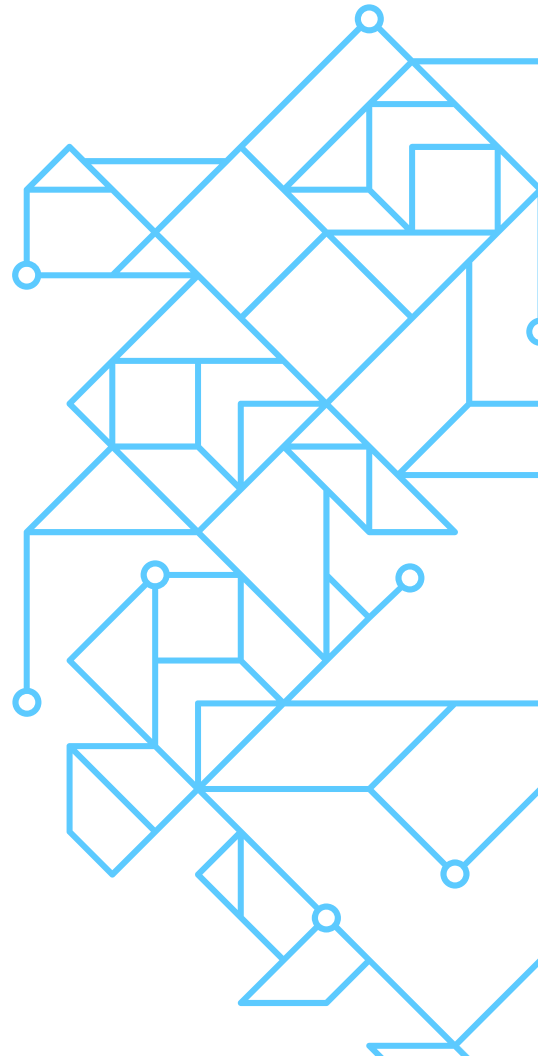**Adam**

adamb@anyscale.com

**Kamil**

kamil@anyscale.com

📋 **The Plan**

**Here's what to expect today.**

📅 # Today's agenda.

- ● What is Ray Serve?
- ● Why use Ray and Ray Serve for scalable AI?
- ● Build complex ML applications with Ray + Serve
- ● Under the hood: features for powering production apps
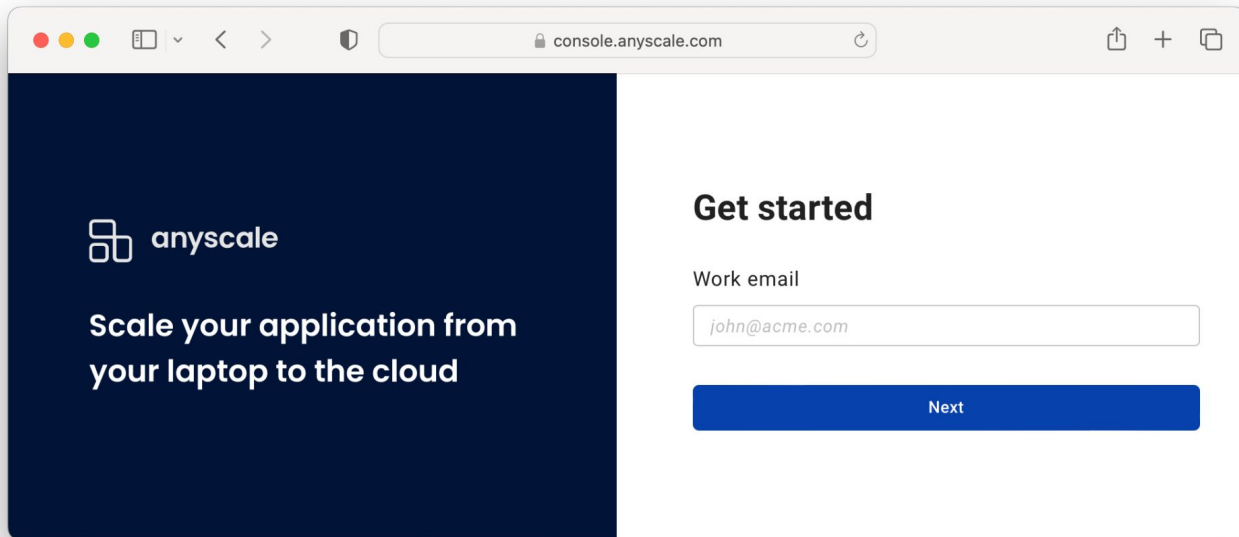- ● Architecture options, hands-on labs, and Q&A

# ✅ Tech check.

## Accessing Anyscale clusters.

- All work will be in Anyscale provisioned clusters.
- Our GitHub repo will be mounted automatically.
- Access begins now.
  - Check your email for login information.
  - Step-by-step instructions to follow.
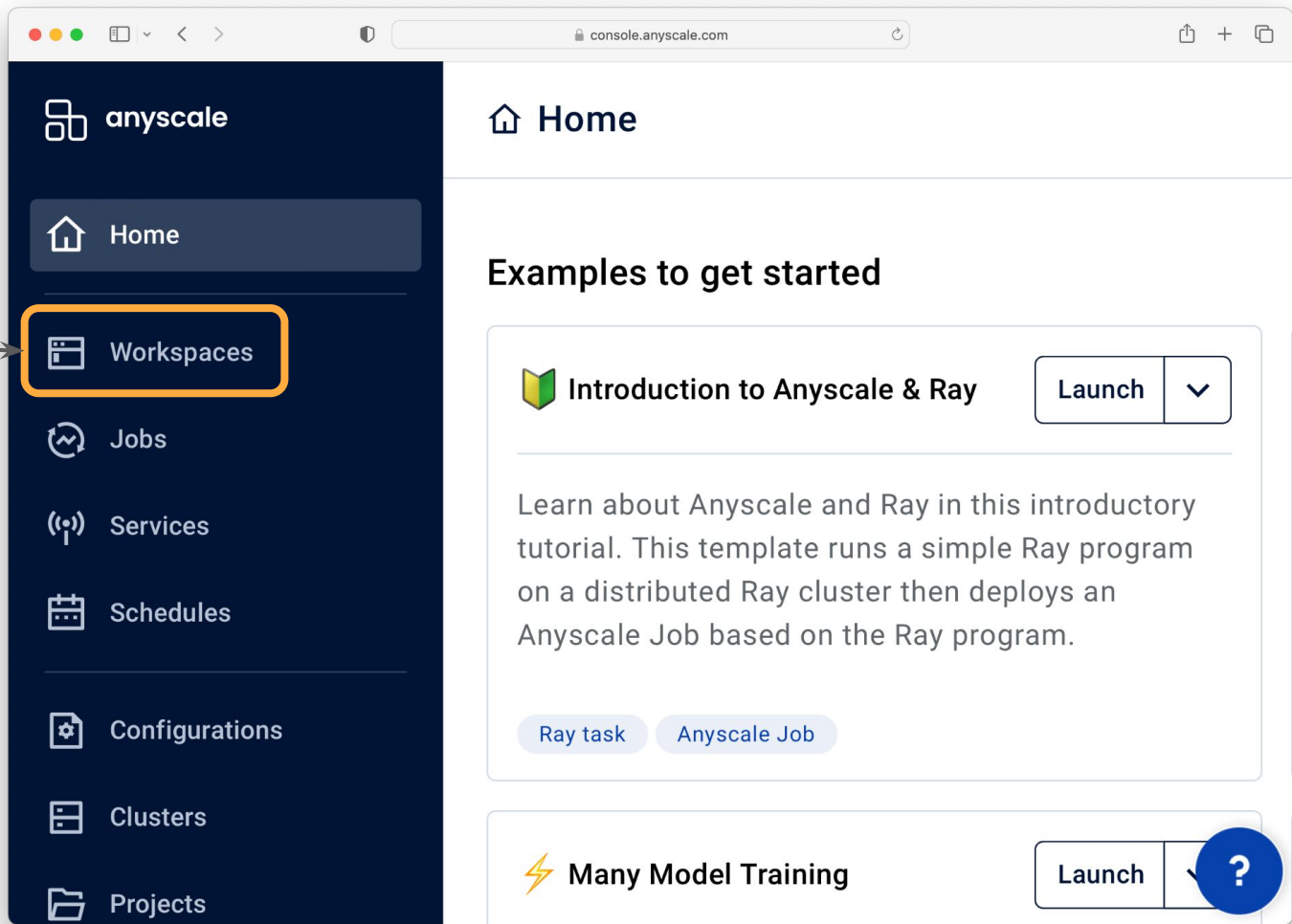
# Anyscale login

Link to Anyscale cluster: console.anyscale.com



Enter the **unique credentials** sent to your email!

1. Select Workspaces

2. Select Your Workspace

3. Click on Jupyter icon

4. Find the content for your class here.

☕ **Time for a Break!**

15 minutes.

🍎 **Today we learned...**

🚂 What Ray Serve is and how it works

🔬 How to use Serve for production services

🔍 Why to choose Serve for AI-based apps

# More Resources

For further exploration with Ray, Anyscale, and LLMs.

# 🔗 Reading list.

### Ray Education GitHub
Access bonus notebooks and scripts about Ray.

### Ray documentation
API references and user guides.

### Anyscale Blogs
Real world use cases and announcements.

### YouTube Tutorials
Video walkthroughs about learning LLMs with Ray.

# 👩‍💻 Connect with the community.

👋 Join the community

*Attend events*, *subscribe to newsletter*, *follow on Twitter*.

🙋 Get support

*Join Ray Slack*, *ask questions on forum*, *open an issue*.

🔭 Contribute to Ray

*Read contributor guide*, *create a pull request*.

# Thank you!

We hope to meet again.