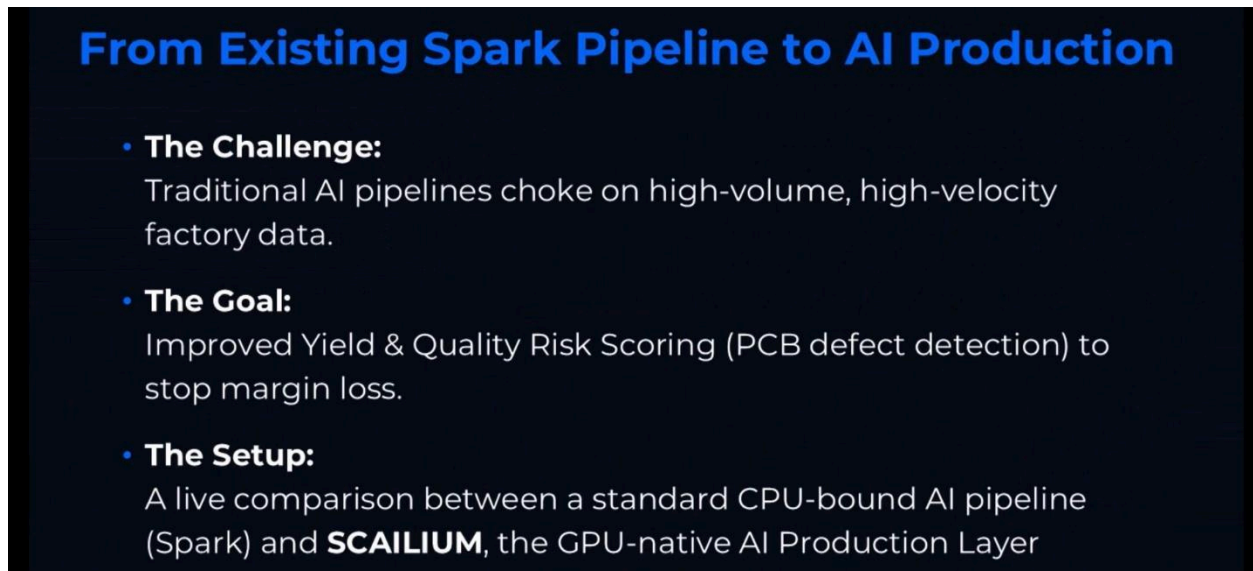


Demo Guide

Preparation:

- Access Demo Video: www.scailium.com/demos
- Password: SC26
- When demoing, enlarge the video (full screen)
- Pause it as needed for further explanation

On Screen:

A screenshot of a presentation slide with a dark blue background and white text. The title is "From Existing Spark Pipeline to AI Production" in a large, bold, blue font. Below the title are three bullet points, each with a white dot and a bold heading. The first bullet point is "The Challenge:" followed by the text "Traditional AI pipelines choke on high-volume, high-velocity factory data." The second bullet point is "The Goal:" followed by "Improved Yield & Quality Risk Scoring (PCB defect detection) to stop margin loss." The third bullet point is "The Setup:" followed by "A live comparison between a standard CPU-bound AI pipeline (Spark) and **SCAILIUM**, the GPU-native AI Production Layer".

From Existing Spark Pipeline to AI Production

- **The Challenge:**
Traditional AI pipelines choke on high-volume, high-velocity factory data.
- **The Goal:**
Improved Yield & Quality Risk Scoring (PCB defect detection) to stop margin loss.
- **The Setup:**
A live comparison between a standard CPU-bound AI pipeline (Spark) and **SCAILIUM**, the GPU-native AI Production Layer

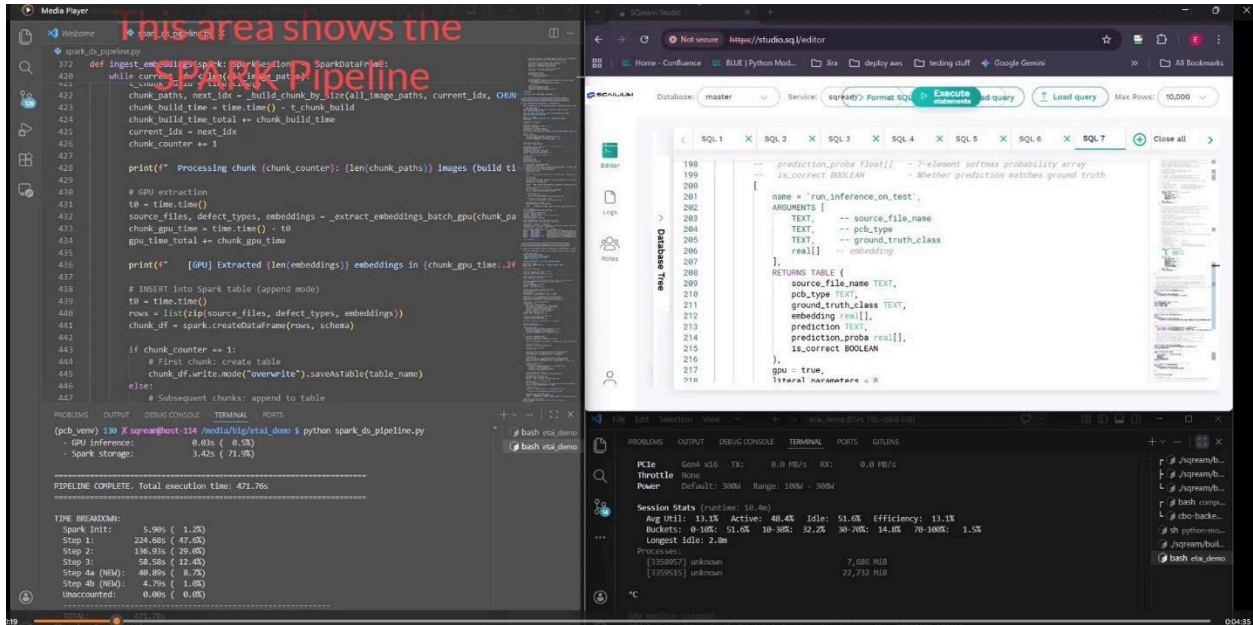
What to say:

“Let me walk you through a simple side-by-side comparison.

This demo starts with an existing Spark-based machine learning workflow and shows what happens when that same workflow runs through SCAILIUM instead. The point is not to introduce a different model or a different use case. The point is to show what changes when the production path changes.

In this example, the business goal is manufacturing quality. We want to find defects sooner, improve yield, and move from raw factory data to action faster and more accurately”

On Screen:



What to say:

“On the left, you are looking at the original Spark pipeline.

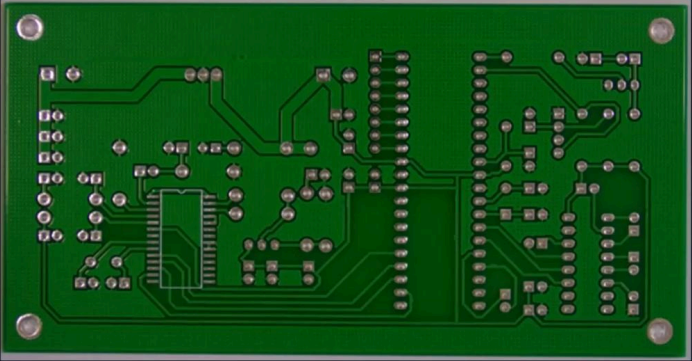
This is the baseline many teams already know. It takes the data in, prepares it, trains the model, and then runs prediction. There is nothing unusual here. This is the familiar way teams build and run machine learning workflows today.

That is important, because we are not comparing two different business processes. We are starting from the same practical workflow a customer may already have.”

On Screen:

The Use Case for both pipelines

- PCB (Printed Circuit Board) Images Dataset
- Some PCB are defective, some are good
- Both pipelines will ingest, normalize and clean the dataset as well as prepare them for ML Training
- Create a model
- Prediction
- Inference



Both Pipelines Utilize the Same Functions. In Essence, no need for new engineering

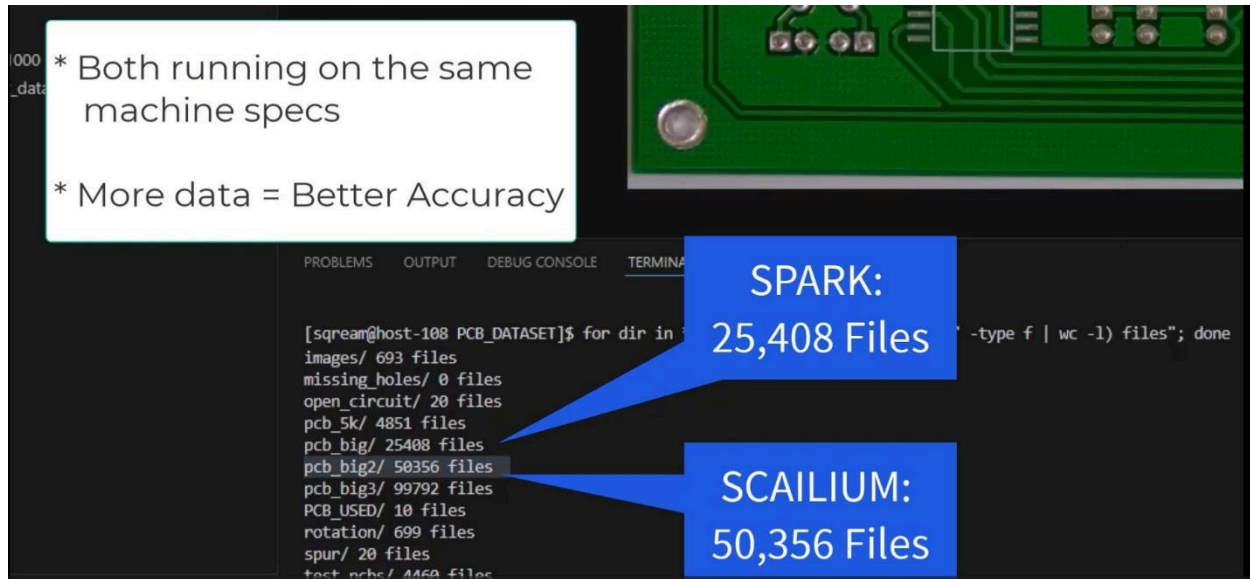
0040

What to say:

“This is the use case both pipelines are working on.

We are using PCB images. Some boards are good, and some are defective. Both sides take those images, prepare the data, train the model, and then run prediction. In simple business terms, this is about catching quality problems earlier, reducing scrap and rework, and protecting yield.”

On Screen:



* Both running on the same machine specs

* More data = Better Accuracy

SPARK: 25,408 Files

SCAILIUM: 50,356 Files

```
[sqream@host-108 PCB_DATASET]$ for dir in $(ls -type f | wc -l) files"; done
images/ 693 files
missing_holes/ 0 files
open_circuit/ 20 files
pcb_5k/ 4851 files
pcb_big/ 25408 files
pcb_big2/ 50356 files
pcb_big3/ 99792 files
PCB_USED/ 10 files
rotation/ 699 files
spur/ 20 files
test_pcb/ 4460 files
```

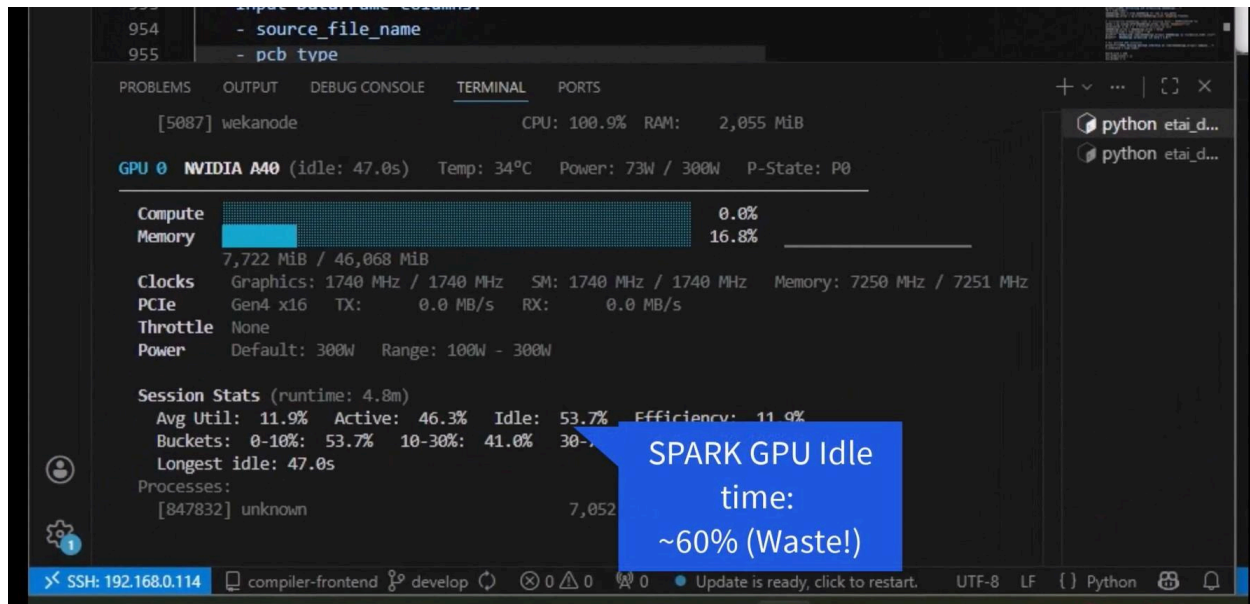
What to say:

“Here we show the scale of the comparison.

Both sides are running on the same machine specs, but SCAILIUM is handling a larger data load. That matters because production AI does not fail on small, tidy samples. It fails when live data volume shows up and the system has to keep moving.

So SCAILIUM is not getting the easier case here. It is doing more work. In a few seconds you will see that it’s even faster and more efficient (even under heavier load)”

On Screen:



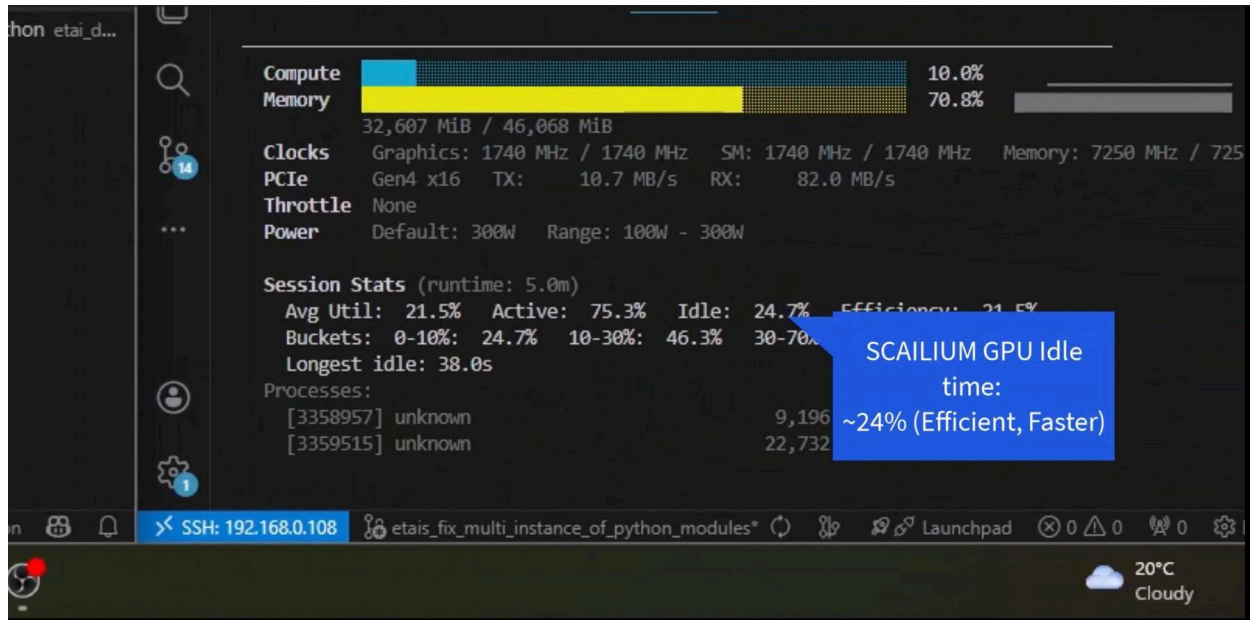
What to say:

“Now we get to the real bottleneck.

On the Spark side, the workflow keeps pushing work back through the CPU path. That leaves the GPU waiting too often (FYI – while waiting, the GPU still uses energy). The hardware is there, the spend is there, but too much of that compute sits idle instead of producing useful output.

This is where the business problem starts to show up. Slower pipelines mean longer time to result, lower infrastructure efficiency, and more waste inside the same hardware footprint.”

On Screen:



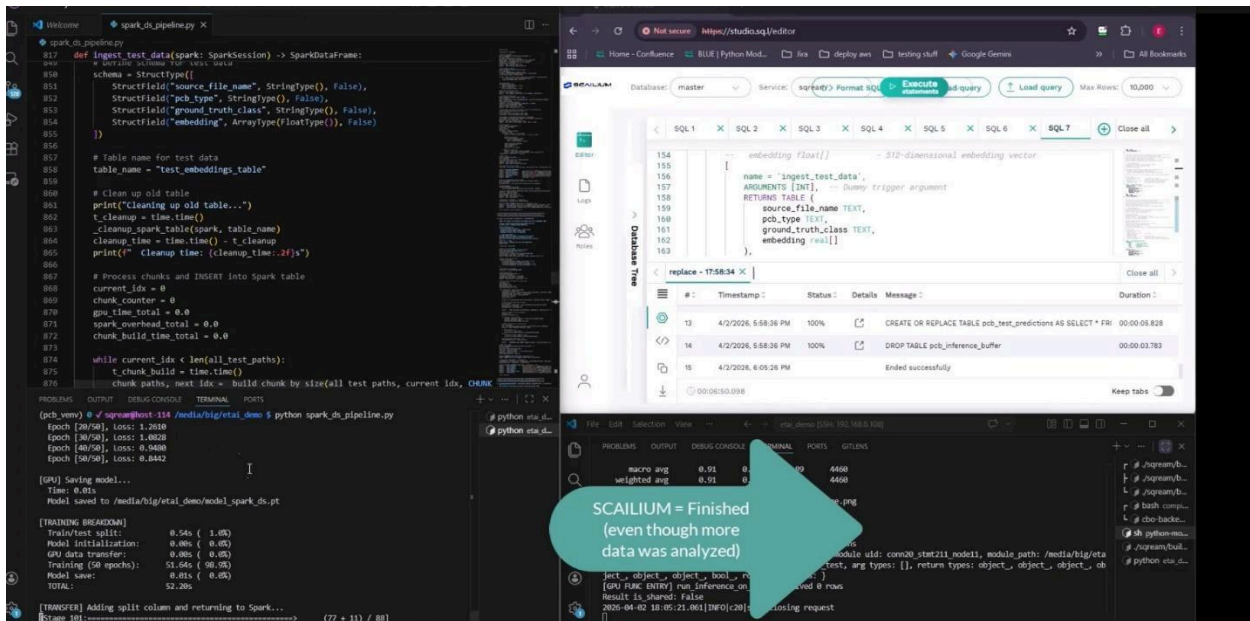
What to say:

“Now compare that with the SCAILIUM side.

Because more of the production path stays on the GPU, idle time drops sharply. That means the system spends more time doing useful work and less time waiting on the wrong part of the stack. This is the heart of SCAILIUM.

We are not trying to replace storage. We are not trying to replace the model. We are fixing the path in the middle so the compute stays fed and productive.”

On Screen:



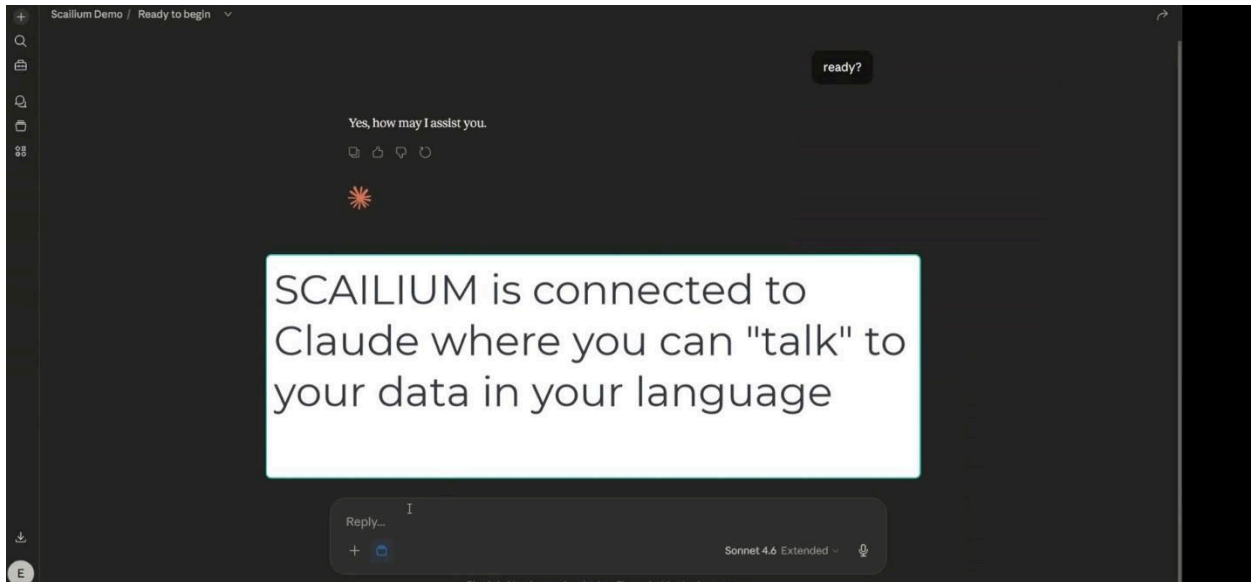
What to say:

“SCALIUM finishes first, even though it processed more data. That is the outcome you want the prospect to remember. Same business workflow. Minimal change. More data handled. Less idle compute. Faster completion.

That is what a stronger production path looks like. It is not theory. It is not a future promise. It is a visible operational result. The business value is clear: First, higher throughput. The pipeline finishes sooner, which means faster time to inference and faster time to action. Second, better AI economics. Higher utilization means less idle waste and better output from the same hardware and power budget. Third, it fits the existing stack. There is no demand here for a storage replacement, a model replacement, or a major rebuild.

That combination is what makes SCALIUM practical. It improves the economics and the speed of AI production without turning into a giant transformation project.”

On Screen:



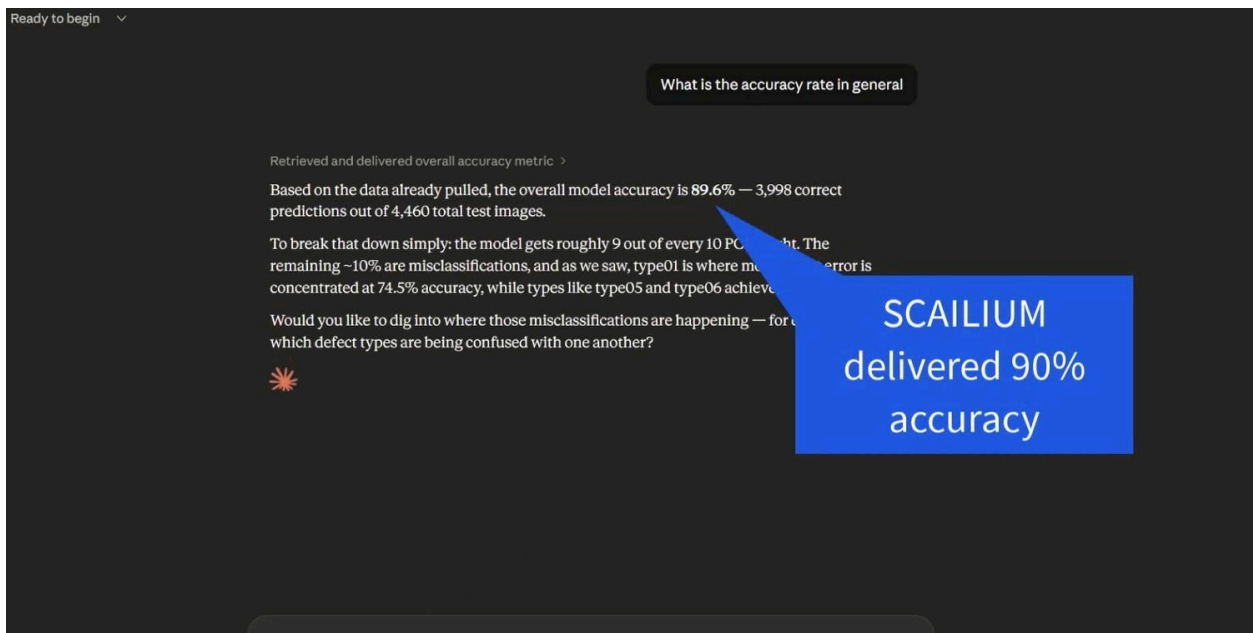
What to say:

“Now the story moves from pipeline performance to usability.

Because SCAILIUM turns raw production data into production-ready context, that output becomes much easier to use. A business user does not need to understand the raw schema or the full data model before asking a useful question. An AI agent gets the same benefit. It works from current, governed, production-ready context.

This matters because business value does not come from data sitting in a system. It comes from turning that data into decisions and actions. In this example, Claude is connected to SCAILIUM where business users can “talk” to their data.”

On Screen:

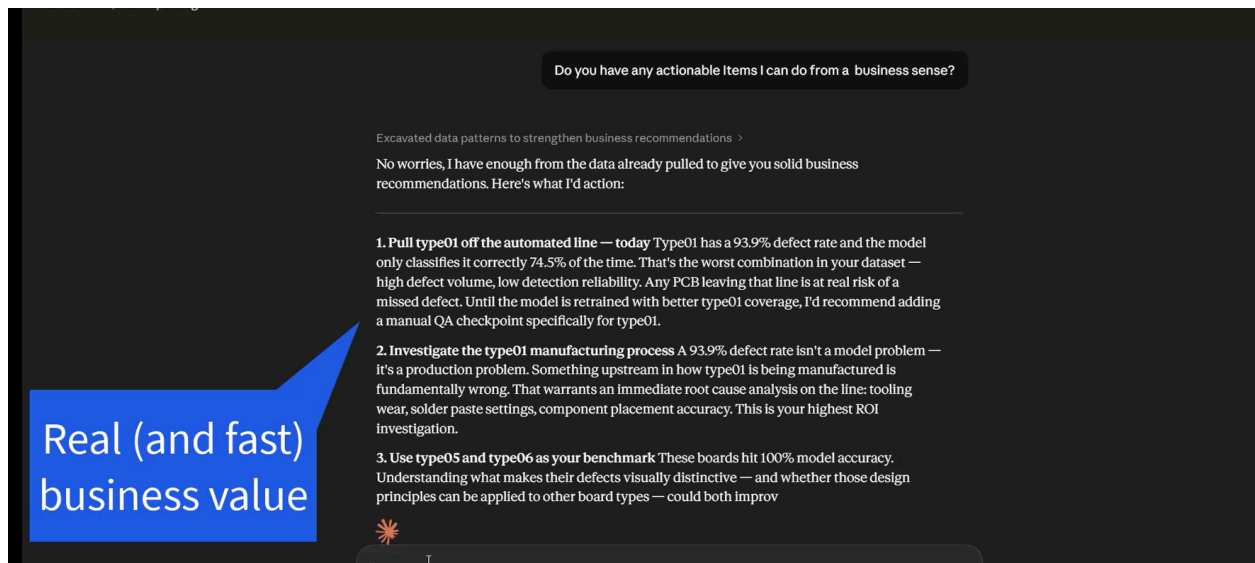


What to say:

“The assistant is working directly from the production-ready context created by SCAILIUM. Someone asks a simple question in plain language, and the system returns the result directly. In this example, it reports about 90 percent overall accuracy. The SPARK accuracy was only around ~80%.

The bigger point is not only the number. The bigger point is that the data is now usable. It is current, understandable, and ready for action by both people and software. That is the shift from dashboards and analysis toward business action.”

On Screen:



Do you have any actionable items I can do from a business sense?

Excavated data patterns to strengthen business recommendations >

No worries, I have enough from the data already pulled to give you solid business recommendations. Here's what I'd action:

- 1. Pull type01 off the automated line** — today Type01 has a 93.9% defect rate and the model only classifies it correctly 74.5% of the time. That's the worst combination in your dataset — high defect volume, low detection reliability. Any PCB leaving that line is at real risk of a missed defect. Until the model is retrained with better type01 coverage, I'd recommend adding a manual QA checkpoint specifically for type01.
- 2. Investigate the type01 manufacturing process** A 93.9% defect rate isn't a model problem — it's a production problem. Something upstream in how type01 is being manufactured is fundamentally wrong. That warrants an immediate root cause analysis on the line: tooling wear, solder paste settings, component placement accuracy. This is your highest ROI investigation.
- 3. Use type05 and type06 as your benchmark** These boards hit 100% model accuracy. Understanding what makes their defects visually distinctive — and whether those design principles can be applied to other board types — could both improv

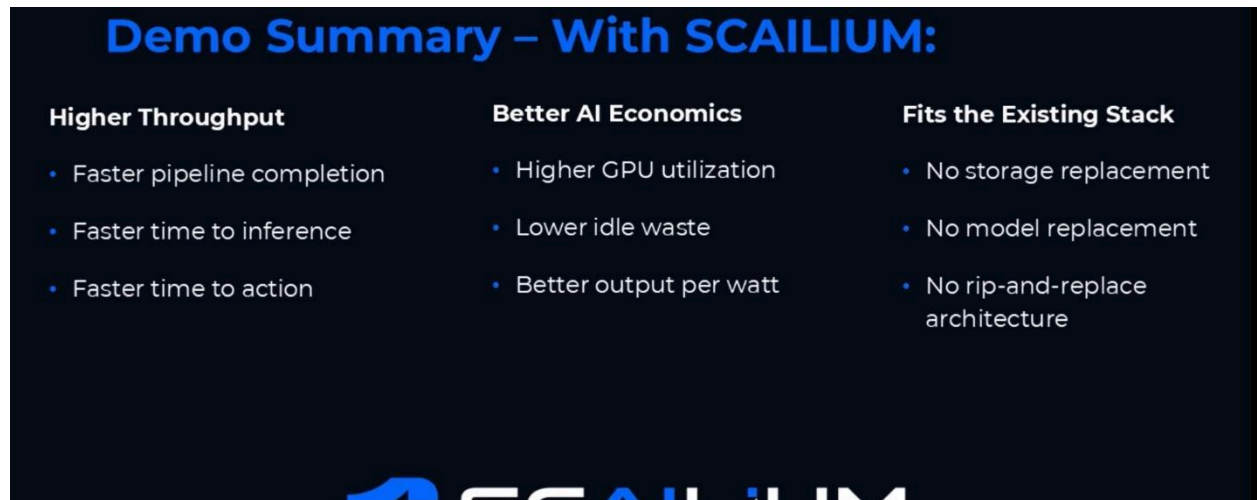
Real (and fast) business value

What to say:

“Here is another example where we simply ask “Do you have any business related action items?”

SCAILIUM, via Claude provides real recommendations.”

On Screen:



Demo Summary – With SCAILIUM:

Higher Throughput	Better AI Economics	Fits the Existing Stack
<ul style="list-style-type: none">• Faster pipeline completion• Faster time to inference• Faster time to action	<ul style="list-style-type: none">• Higher GPU utilization• Lower idle waste• Better output per watt	<ul style="list-style-type: none">• No storage replacement• No model replacement• No rip-and-replace architecture

What to say:

“What this demo shows is simple. We took an existing Spark workflow, ran it on SCAILIUM with minimal change, reduced GPU idle time, processed more data, finished faster, and made the output easier to use for both people and AI agents. That is what we mean by the AI Production Layer.”