**EPISODE #346**

# GPU Scaling: The "Go Wide or Go Tall" Dilemma

Published January 29, 2026 • Runtime: 25:19

https://myweirdprompts.com/episode/serverless-gpu-scaling-efficiency/

## EPISODE SYNOPSIS

In this episode, Herman and Corn dive deep into the engineering trade-offs of serverless GPU workloads. Using a real-world text-to-speech example on the Modal platform, they explore whether it's better to scale horizontally with many small workers or vertically with a single high-end GPU like the H100. They break down the hidden costs of cold starts, the importance of memory bandwidth over raw compute, and how to find the "sweet spot" on the cost-efficiency curve to get the most bang for your buck.

## DANIEL'S PROMPT

**Daniel**

I'm looking for advice on cost optimization for serverless GPU functions. What is more cost-effective: parallelizing tasks by spawning multiple workers to run simultaneously or using a more powerful GPU with more VRAM to shorten generation time, despite the higher hourly cost? I'm trying to figure out how to get the maximum bang for my serverless budget.

# TRANSCRIPT

### Corn

Herman, I think our housemate Daniel is trying to either bankrupt us with his curiosity or save us a fortune. He sent over a prompt that gets right into the weeds of how this very show is put together.

### Herman

Herman Poppleberry at your service. And yeah, I saw that one come in. It is a classic engineering trade-off, isn't it? Daniel is basically asking the million-dollar question for anyone building in the AI space right now. When you are running these heavy workloads on serverless GPUs, do you go wide or do you go tall?

### Corn

Right, and specifically for our listeners who might not be wrestling with serverless configurations every morning, we are talking about how to handle tasks like generating the audio for this podcast. Daniel mentioned he is using Modal, which is a great serverless platform, to run these text to speech functions. The dilemma is: do you spawn ten small, cheap workers to do ten things at once, or do you pay for one absolute monster of a GPU to finish the whole job in a fraction of the time?

### Herman

It is the Ferrari versus a fleet of scooters problem. Do you want one package delivered at two hundred miles per hour, or fifty packages delivered at thirty miles per hour? And in the serverless world, where you are billed by the second, the math is not always as intuitive as it seems.

### Corn

Well, let's break that math down because I think that is where the "bang for the buck" really lives. Most people assume that if a GPU is four times as expensive per hour, it must be four times as fast to break even. But it is rarely a linear relationship, is it?

**Herman**

Exactly. Let's look at the hardware for a second. On platforms like Modal in early twenty twenty-six, you usually have a choice between something like an Nvidia L four on the lower end, and then you jump up to the H one hundred or even the new Blackwell B two hundred on the high end. An L four might cost you about eighty cents an hour, while a B two hundred with its massive memory is going to be significantly more—closer to six dollars an hour. But here is the thing most people don't realize: VRAM, or video random access memory, isn't just about capacity. It is about bandwidth.

**Corn**

Right, so it is not just how much data the GPU can hold, but how fast it can shove that data into the processing cores. If you are running a large language model or a complex text to speech engine, that bandwidth is often your primary bottleneck, not the raw compute power.

**Herman**

Spot on. If your model is constantly waiting for data to move from memory to the chip, you are paying for the GPU to sit idle for milliseconds at a time. On an H two hundred, that bandwidth is massive—nearly four point eight terabytes per second. On an older L four, you are looking at a fraction of that, maybe three hundred gigabytes per second. So, even though the H two hundred costs more per second, it might finish the task fifteen times faster because it is actually staying busy.

**Corn**

But wait, Herman, Daniel's question also brought up parallelization. If he spawns, say, twenty workers on cheaper L fours, they all start at the same time. In theory, the total wall-clock time for the whole project drops significantly. Does the overhead of starting those twenty containers eat up the cost savings?

**Herman**

That is the "cold start" trap. In a serverless environment, when you tell the provider to "spawn a worker," it has to find a machine, pull your container image, load your model into memory, and then start processing. Even with Modal's new GPU snapshotting, which has cut cold starts down to about ten seconds, if your task only takes two seconds of actual compute, you just paid for twelve seconds of time for two seconds of work. If you do that twenty times in parallel, you are paying for two hundred and forty seconds of "setup" time.

## Corn

Okay, so if the task itself is very short, parallelizing across many small workers is actually a massive waste of money because you are multiplying the cost of the cold starts.

## Herman

Precisely. This is what we call the "granularity of the task." If Daniel is processing this podcast script sentence by sentence, and each sentence takes two seconds to generate, he should absolutely not be spawning a new worker for every sentence. He would be spending eighty percent of his budget on the overhead of spinning up the environment.

## Corn

That makes a lot of sense. So, there is a sweet spot where the task is long enough that the cold start becomes negligible, but short enough that you still benefit from parallelism. I am curious about the VRAM aspect too. Daniel mentioned "more VRAM to shorten generation time." Does having "extra" VRAM help if the model already fits?

## Herman

That is a common misconception. If your model needs twelve gigabytes of VRAM and you give it sixteen, it won't necessarily run faster than if you gave it eighty, unless that eighty-gigabyte card also has higher memory bandwidth or more CUDA cores. However, more VRAM allows for larger batch sizes. Instead of generating one voice line at a time, you could potentially generate sixteen lines simultaneously on one big GPU like an H one hundred.

## Corn

Ah, so "vertical scaling" where you use a bigger card actually enables a different kind of "internal parallelism" through batching.

**Herman**

Exactly. And batching is almost always more cost-effective than spawning separate workers. When you batch, you only pay for one cold start, one model load, and one idle period. The GPU cores are designed to handle multiple operations at once. If you are only using ten percent of a B two hundred's capacity to generate one line, you are wasting ninety percent of what you are paying for. But if you can saturate that card with a large batch, your cost per generated second of audio drops through the floor.

**Corn**

So, for Daniel's specific use case with the podcast audio, the strategy should probably be: grouping the lines into chunks, using a powerful enough GPU to handle those chunks in batches, and keeping the worker "warm" if possible.

**Herman**

Right. Modal has this great feature where you can set a "container idle timeout" or use "request batching." If you know you are going to be sending more lines in the next few minutes, you keep the container alive. You pay for the idle time, yes, but you avoid the ten-second cold start on the next request. It is all about the ratio of compute time to overhead time.

**Corn**

I remember back in episode three hundred and seventeen, we talked about the sunk cost trap. In computing, people often get stuck in the "sunk overhead trap." They've already spent the time and money to get the model loaded, so they feel like they should just keep using that one worker even if it's slow. But sometimes, switching to a more expensive, faster instance actually clears the queue so fast that the total bill is lower.

**Herman**

It is counter-intuitive until you see the logs. I've seen cases where moving from an L four to an H one hundred—which is a much more powerful card—cut the execution time by eighty percent while only increasing the hourly rate by four times. That is a net win for the budget.

**Corn**

Let's talk about the "bang for the buck" benchmarks. If someone is listening and they are trying to optimize their own serverless bill, what is the first metric they should look at? Is it "seconds per inference" or "dollars per inference"?

**Herman**

It has to be dollars per inference. You take the hourly rate of the GPU, divide it by three thousand six hundred to get the cost per second, and then multiply that by how many seconds the task took, including the setup time. Do that for a small GPU, a medium GPU, and a large GPU. You will often find a "U-shaped" curve.

**Corn**

A U-shaped curve? Explain that.

**Herman**

Well, on the very cheap GPUs, the cost per inference is high because the task takes forever. On the very expensive GPUs, the cost per inference might also be high because you aren't fully utilizing the massive power of the card—you're paying for "empty seats" in the stadium, so to speak. The "bottom" of that U-curve is your sweet spot. Usually, for something like text to speech or image generation in twenty twenty-six, it is a mid-tier card like the L four or a partitioned A one hundred if your provider allows for fractional GPUs.

**Corn**

That is a great analogy. You don't want to rent a whole stadium for a ten-person meeting, but you also don't want to hold a thousand-person conference in a coffee shop where everyone has to wait in line for an hour to get a seat.

**Herman**

Exactly. And there is another factor we haven't touched on: data transfer. In serverless functions, you are often moving large model files from a storage bucket into the GPU's local memory. If you spawn fifty workers, you are potentially being billed for the time it takes to download those fifty copies of the model.

**Corn**

Oh, that is a huge point. If the model is five gigabytes, and the network speed is limited, you could be spending more time watching a progress bar than actually generating audio.

**Herman**

And on some platforms, you pay for that network egress or ingress too. So, the "fleet of scooters" approach starts looking very expensive when every scooter has to stop at the gas station and fill up its own five-gallon tank before it can start the delivery. One big truck with one big tank is more efficient for heavy loads.

**Corn**

So, to summarize the advice for Daniel and anyone in his shoes: if the total workload is large, you should aim for the most powerful GPU you can reasonably saturate with batching. Only use parallel workers if the tasks are truly independent, long-running, and you've already optimized the cold start.

**Herman**

Right. And for text to speech specifically, which is what Daniel is doing, the latency often matters for the user experience. If he is generating these episodes and he wants to hear the result quickly, the "wall-clock time" matters. Parallelization wins on wall-clock time, but vertical scaling often wins on cost-efficiency.

**Corn**

It is the classic "Time, Cost, Quality" triangle, though in this case, the quality is usually the same regardless of the GPU, so it is really just a trade-off between how fast you want it and how much you want to pay.

**Herman**

Actually, Corn, I would argue that in some cases, the more powerful GPU can lead to better quality if it allows you to use a larger, more sophisticated model that simply wouldn't fit or would be too slow on a smaller card. If you are restricted to a twenty-four gigabyte L four, you might have to use a "distilled" or "quantized" version of a model. If you step up to a one hundred and forty-one gigabyte H two hundred, you can run the full, uncompressed weights.

**Corn**

That is a fair point. We often forget that hardware limits the software we even consider using. If Daniel knows he has the budget for a B two hundred, he might choose a much more advanced audio synthesis model that sounds more human.

**Herman**

Exactly. It opens up the "possibility space." But back to the "bang for the buck" for the current setup—I genuinely think he should look into "request batching." On Modal, you can actually set up a function to wait a few milliseconds to see if more requests come in, and then process them all together. It is like waiting for the elevator to fill up before it goes to the top floor.

**Corn**

I love that. It is much better than the elevator going up and down for every single person. So, Herman, if you had to give Daniel a "rule of thumb" for his next experiment, what would it be?

**Herman**

I would say: aim for seventy percent utilization. If your GPU is sitting at ten percent load while it generates a single line of dialogue, you are over-provisioned. Increase your batch size or move to a smaller, cheaper card. If your card is at ninety-nine percent and the task is taking ten minutes, try the next tier up. You might find that the total cost actually stays the same or drops because the time-saving is so significant.

**Corn**

It is about finding that "flow state" for the hardware. I think people get intimidated by the hourly rates of high-end GPUs. Seeing "six dollars an hour" looks scary compared to "eighty cents an hour." But if the six-dollar card finishes in thirty seconds and the eighty-cent card takes twenty minutes... well, you do the math.

**Herman**

Exactly. Twenty minutes at eighty cents an hour is about twenty-seven cents. Thirty seconds at six dollars an hour is about five cents. The "expensive" card is actually five times cheaper in that scenario.

**Corn**

That is the "insider secret" right there. The sticker price is not the cost of the job. It is just the rate of the meter. If the car is ten times faster, a higher meter rate can still mean a cheaper fare.

**Herman**

And that is why we do this show, Corn. To make sure people aren't overpaying for their metaphorical taxi rides in the cloud.

**Corn**

Well, speaking of the show, I think it is fascinating that we are talking about the very infrastructure that makes our voices possible. It is a bit meta. Daniel is out there tweaking the knobs on the machine that lets us exist in this digital space.

**Herman**

It makes me feel a bit like a high-maintenance digital entity. "Excuse me, Daniel, I require a B two hundred today. My thoughts are far too complex for an L four."

**Corn**

You've always been a Blackwell in my book, Herman. High bandwidth, lots of memory, occasionally runs a bit hot.

**Herman**

I will take that as a compliment. But seriously, the development in this space just in the last year is staggering. We are seeing these "serverless" providers get much better at reducing those cold start times. Modal's snapshotting is a game changer because it saves the entire state of the GPU memory, so the "setup" time is effectively just the time it takes to resume a process.

**Corn**

If the setup time goes to zero, does that change the math back in favor of parallelization?

**Herman**

It does! If there is no penalty for spawning a new worker, then you should always parallelize to the maximum extent possible to get the fastest result. The only reason we don't do that now is because of the "tax" we pay every time a new container starts.

**Corn**

So, the future of cost optimization is actually the providers getting more efficient at the "boring" stuff like container orchestration and image pulling.

**Herman**

Exactly. We are moving toward a world where compute is truly a utility, like turning on a light switch. You don't pay for the power plant to warm up; you just pay for the electrons you use the moment you flip the switch. We aren't quite there yet with GPUs because they are such specialized, scarce resources, but we are getting closer.

**Corn**

It reminds me of episode one hundred and twenty-eight, where we looked back from the future at the "Dial-Up Era" of AI. Right now, waiting ten seconds for a GPU to start feels like waiting for a modem to hand-shake. In ten years, we will laugh at the idea that we had to "wait" for a model to load.

**Herman**

I can't wait for that. I want my insights to be instant. But until then, Daniel has to be the architect. He has to balance the budget against the clock.

**Corn**

I think he is doing a pretty good job. This episode sounds great so far, at least from where I am sitting.

**Herman**

Agreed. And for anyone else out there wrestling with these serverless prompts, the best advice is to log everything. Don't guess. Modal gives you those beautiful logs that show exactly when the container started, when the model finished loading, and when the inference ended. Use that data. It is your best friend.

**Corn**

Data-driven optimization. It is not as flashy as "AI," but it is what actually builds sustainable businesses and projects.

**Herman**

And it keeps the lights on here at the Poppleberry household in Jerusalem.

**Corn**

Exactly. Well, Herman, I think we've given Daniel plenty to chew on for his next round of optimization. I'm curious to see if the next episode arrives in my inbox five minutes faster.

**Herman**

If he switches to B two hundreds, we might finish the next episode before he even sends the prompt. That is the dream.

**Corn**

Temporal anomalies aside, I think we should wrap this one up. But before we go, I want to say a huge thank you to everyone who has been following along with "My Weird Prompts." We've hit episode three hundred and forty, which is just wild to think about.

**Herman**

It really is. We've covered everything from ADHD to North Korean tech to WORM drives, and it is all thanks to the weird and wonderful questions Daniel sends our way. And of course, the fact that you all keep listening.

**Corn**

If you are enjoying the deep dives, we would really appreciate it if you could leave us a review on Spotify or wherever you get your podcasts. It genuinely helps the show grow and helps other curious people find these rabbit holes we love to explore.

**Herman**

Yeah, a quick rating or a few words about your favorite episode makes a big difference. We read them all, and it keeps us motivated to keep digging into these technical nuances.

**Corn**

You can find all our past episodes and a way to get in touch with us at our website, myweirdprompts.com. We love hearing from listeners, even if you just want to tell us your own GPU optimization horror stories.

**Herman**

Or if you have a prompt that is even weirder than Daniel's. We are always up for a challenge.

**Corn**

Alright, that is it for today. Thanks for joining us for another episode of My Weird Prompts. I am Corn.

**Herman**

And I am Herman Poppleberry. We will see you in the next one.

**Corn**

Stay curious, everyone.

**Herman**

And keep those GPUs saturated.

**Corn**

See ya!

**Herman**

Bye!