



AI

RESEARCH ON NVIDIA DGX SPARK

Pushing the Frontier of Local AI
on a Petascale Desktop

Manav Sehgal

Contents

1	AI Research on NVIDIA DGX Spark	5
2	Derisking the Cloud Pretrain — How a \$5K Spark Saves \$50K on H100 Rentals	7
2.1	Why this matters for a personal AI builder	7
2.2	Where the cloud cost actually lives	8
2.3	Phase 1 — the Spark recipe lab (~\$1, ~2 hours)	9
2.4	Phase 2 — the Spark stress test (~\$0.40, ~22 hours)	10
2.5	Phase 3 — the cloud commit (~\$3,360, ~1 week)	10
2.6	Verification — what success looks like on the Spark for a cloud-bound workload	11
2.7	Tradeoffs and gotchas	12
2.8	What this unlocks	13
2.9	State of the apps — as of this article	13
3	Keep reading	15

1. AI Research on NVIDIA DGX Spark

Pushing the Frontier of Local AI on a Petascale Desktop

by Manav Sehgal

This is a free sample from the full book.

2. Derisking the Cloud Pretrain — How a \$5K Spark Saves \$50K on H100 Rentals

The Spark is too small for a serious pretrain — but it's the right size for the recipe-search that precedes one. Cull 100 candidate architectures down to 3 on one Spark for ~\$1 of electricity, then book the cloud node knowing what to train. The expected savings per campaign run into the thousands.

You cannot pretrain a 7-billion-parameter Llama-class model on a DGX Spark from scratch. A Chinchilla-optimal 7B run needs around 140 billion training tokens — six times the parameter count, ten times if you intend to overtrain — and the FLOP budget that comes with it. On one GB10 the wall clock is months; on a single H100 it is weeks; on a properly-spec'd 8-GPU cloud node it is days. The cloud is the right tool for that workload. The Spark is not.

The argument of this article is that this fact about the Spark is *not* its weakness. It is the precondition for the Spark's most undersold use case — the **recipe lab** that decides what the cloud node should actually run when you finally book it. A pretraining campaign is not one job; it is a hundred candidate jobs collapsing to one. The cloud is the right tool for the *one*. The Spark is the right tool for the *hundred*.

The arithmetic that follows is the case for treating your DGX Spark as a wind tunnel before any aircraft is built. A hundred-iteration architectural sweep on one Spark costs about a dollar of marginal electricity and a couple of hours of wall time. The same sweep at meaningful scale on cloud — single-H100 spot, sustained — runs ~\$5 to ~\$70 depending on how target-faithful you make it. The much larger number is what you save on the *back end*: the wrong-architecture cloud campaign you didn't book because the Spark already told you it would not work. At a 50% wrong-pick rate without prior signal — conservative for blind architectural search — a \$1 Spark sweep gates ~\$1,679 in expected loss against a small (\$3K) cloud campaign and ~\$7K against a medium (\$14K) one. Scaled to a 70B Llama-class run on 1024 H100s for 21 days, the same Spark dollar gates more than a million.

📖 Pretraining campaign

A single end-to-end run that trains a base language model from random initialization to a target loss on a fixed token budget. *Campaign* (rather than *job*) emphasizes that the work is one decision-point with a single bill — typically days to weeks of multi-GPU wall-clock and a four-to-seven-figure cloud invoice. Distinct from fine-tuning (which adapts an existing base) and from continued pretraining (which extends an existing base on new data).

📖 Chinchilla-optimal token budget

The DeepMind scaling-laws result (Hoffmann et al., 2022): for a fixed compute budget, the loss-minimizing model trains on ~20 tokens per parameter — roughly 6× more data than the prior Kaplan-laws prescription. At 7B params, “Chinchilla-optimal” means ~140B training tokens; at 70B, ~1.4T. The number sets the lower bound on Phase 3's wall-clock and is the reason the cloud bill exists.

2.1 Why this matters for a personal AI builder

The audience for this article is the engineer at a startup, the research lead at a small lab, or the independent builder who has cloud credits or budget — between \$5K and \$500K — earmarked for *one* serious pretraining run, and exactly one chance to spend it well. That decision is downstream of an architectural search the cloud is a poor place to perform. Doing it on the cloud means paying \$50 to \$5,000 *per candidate* to discover

that your candidate doesn't work. Doing it on the Spark means paying ~\$0.01 per candidate to discover the same thing. The architectural search isn't where the money is — but it's where the *information* is, and information is what compounds.

The Spark's uber move for a personal AI power user has, until this article, been about prototyping at a scale your desk can afford (*fine-tuning a 100B Nemotron*, sized being the prior chapter). Derisking a cloud pretrain is the same idea applied one level up: the Spark on your desk is the rig that lets you walk into the cloud-vendor billing portal with a recipe, a measured trajectory, and a defensible architectural argument — instead of a prayer.

2.2 Where the cloud cost actually lives

The cloud bill for a serious pretrain decomposes into three pieces, and only one of them is what most people think it is.

1. **The final-training run itself.** The headline number — $8 \times \text{H100} \times 168 \text{ hr} \times \$2.50/\text{hr} \approx \$3.4\text{K}$ for a 7B-class spot run, several \times that on-demand or for larger models. This is what gets quoted when someone says “we trained X for \$Y.”
2. **The architectural search that preceded it.** Usually buried in the same line item. If you ran 30 candidate architectures at any meaningful scale before settling, that's another $30 \times \$30$ to $30 \times \$300$ depending on how target-faithful the search was. People don't quote this number because they wish they hadn't spent it.
3. **The expected loss from picking wrong.** The hidden one. If the architectural search wasn't informative — or didn't happen at all — the wrong-pick rate is well above zero, and the cost is the *next* final-training run after the first one didn't converge. At a 50% wrong-pick rate, the expected loss is half of one final-training run per campaign. Anyone who has watched a colleague restart a \$30K cloud booking three weeks in knows this number is not abstract.

The Spark eliminates the second of those three line items entirely (replaces \$50–\$5K of cloud sweep with ~\$1 of Spark electricity) and reduces the *expected value* of the third by an integer factor (cuts the wrong-pick rate from “blind” to “informed”). The first line item it does not touch — you still pay the cloud bill for the final-training run. But you pay it once, on the right architecture, with measured early-loss curves to back the decision.

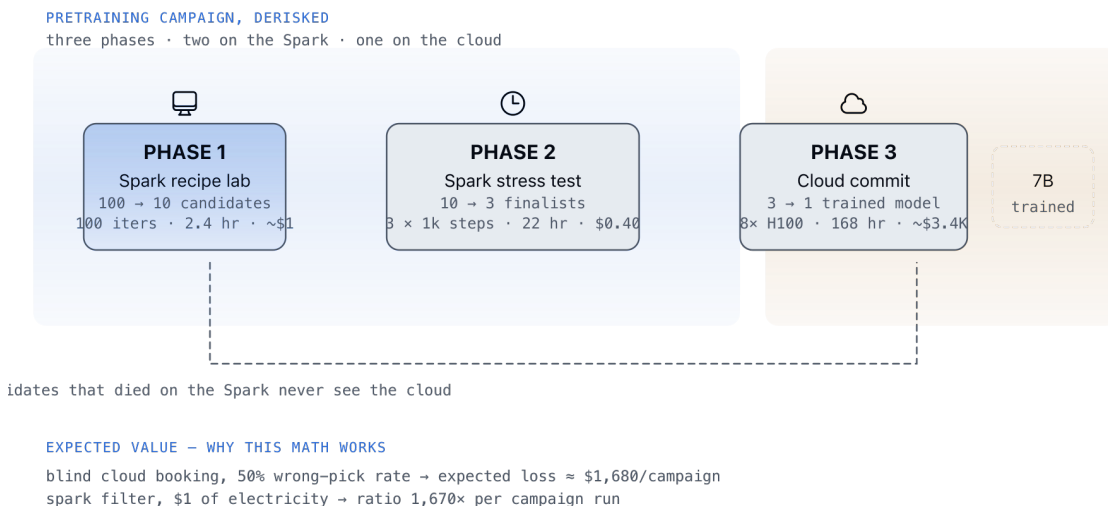


Figure 2.1: A three-phase timeline showing how a cloud pretraining campaign decomposes when the Spark is part of the toolchain. Phase 1 — Spark recipe lab — runs 100 candidate architectures via the A4 agent loop, 88 seconds per iteration, 2.4 hours total wall, costing about one dollar of marginal electricity, culling to 10 promising candidates. Phase 2 — Spark stress test — runs the top 3 candidates for 1000 steps each, 22 hours wall, costing 40 cents, culling to 3 finalists. Phase 3 — cloud commit — books an 8×H100 node for one week to train the chosen architecture to Chinchilla optimal, costing 3,360 dollars at spot rates. A reject path under Phase 1 and Phase 2 shows that candidates that died on the Spark never see the cloud.

2.3 Phase 1 — the Spark recipe lab (~\$1, ~2 hours)

Phase 1 is the **autoresearch agent loop** pointed at a 7B-shaped target architecture instead of an open exploration. The agent proposes a candidate config, the trainer runs it for 60 steps on a proxy model — small enough that one iteration completes in roughly 88 seconds on a Spark — and the loop records `val_bpb` against the trajectory. After 100 iterations, the trajectory log holds 100 (candidate, `val_bpb`) pairs. Sort, take the top 10, save the rest.

What makes Phase 1 cheap is not just the Spark’s electricity cost. It is the **proxy substitution** that the agent loop already encodes: the search runs on a 200M-class model whose head dimension, FFN ratio, and activation function all match the 7B target, but whose total parameter count is two orders of magnitude smaller. The shape of the architectural search (which knob settings dominate, where the loss surface is convex, which combinations interact) transfers up to the target. The absolute numbers do not need to. This is the same trick that makes it possible to run scaling-law experiments at 50M and project to 70B — it is what every serious pretrain campaign already does on the cloud, just done on a desk for 1/300th the cost.

Proxy substitution

Train a smaller model (50M–500M params) that shares the *shape* of the target — same depth-to-width ratio, same activation, same attention pattern — but a fraction of the parameter count. Architectural rankings (which combo of knobs converges fastest) are largely shape-invariant; absolute loss values are not. The proxy tells you which recipe to commit to; the cloud run gets the absolute numbers. The trick is what makes scaling-laws research economically possible.

Why \$1 of Spark electricity gates >\$1,000 of expected cloud savings

The cost being measured isn’t the Spark’s wall draw — it’s the cloud campaign the Spark filtered. A blind cloud booking has a real wrong-pick rate (50% is conservative for architectural search without prior signal). Half of one wrong campaign is the *expected loss*; the Spark’s \$1 cost gates that loss away. The savings ratio doesn’t depend on the Spark being fast — it depends on the cloud being expensive. Frame the Spark as a \$1 lottery ticket whose payoff is “we don’t book the wrong \$14K campaign,” and the math always works in its favor.

The cost arithmetic for Phase 1 lives in `evidence/cost_arithmetic.py` and runs from environment variables so it stays current as cloud spot prices move. The defaults — H100 at \$2.50/hr per GPU, Spark at 240 W sustained training draw, \$0.13/kWh, 1.5-year amortization horizon — produce the table below. The column the article is selling is the Spark column. The rightmost column is the alternative the Spark replaces.

```
$ python3 evidence/cost_arithmetic.py | jq .campaign_100_iters_usd
{
  "spark_electricity_only":      0.08,
  "spark_total_cost_of_use":    1.01,
  "h100_single_proxy":         6.11,
  "h100_8gpu_node_proxy":      48.89,
  "h100_8gpu_node_target_scale": 2933.33
}
```

The Spark electricity-only line (\$0.08 for 100 iterations) is the marginal cost — the Spark is on your desk anyway. The total cost-of-use line (\$1.01) includes amortized hardware over a 1.5-year horizon and is the more honest “what did this run actually cost.” The cloud lines escalate: \$6 to do the same 100-iter sweep on a single rented H100 at proxy scale, \$49 if you parallelize across an 8-GPU node for speed, \$2,933 if you run the same sweep at *target* scale on the cloud (60× longer per iter because the model is 60× larger). The savings ratio depends on which alternative you score against; either way, the Spark column is where you start.

2.4 Phase 2 — the Spark stress test (~\$0.40, ~22 hours)

Phase 1's 60-step taste test catches gross failures and ranks survivors by early-loss slope. It does not catch issues that only emerge at 1,000+ steps — late-onset divergence, optimizer instabilities, gradient norm blowups under longer schedules. Phase 2 is the unattended stress test: run each of the top 3 candidates for 1,000 steps on the Spark and watch for the failure modes that take time to show up.

22 hours of Spark wall clock is one overnight. At the same total cost-of-use rate as Phase 1, that is roughly \$0.40 — call it forty cents to confirm the recipe is stable past the kernel-timing window. If two of the three candidates survive, you have a primary plus a fallback for the cloud booking. If all three survive, you ship the highest-ranked one and keep the others as restart candidates. If only one survives, you have learned something genuinely useful before booking a cloud node: that the architectural envelope is narrow and the next sweep should explore tighter neighborhoods.

The reason Phase 2 belongs on the Spark and not on the cloud is the same reason as Phase 1, just at a longer time horizon. A 1,000-step run on a single cloud H100 costs ~\$0.30 per candidate at proxy scale (close to break-even with the Spark) and ~\$18 at target scale on an 8-GPU node. The cloud is competitive on cost for Phase 2 only if you don't have a Spark. Once you do, the math says: spend the watt-hours, not the dollars.

2.5 Phase 3 — the cloud commit (~\$3,360, ~1 week)

Phase 3 is the only piece that the Spark cannot do — and it is the only piece you should pay full cloud price for. Book the 8× H100 node, run the surviving best candidate to Chinchilla-optimal (or whatever schedule you've decided on), and watch the loss curve land where the Spark stress test predicted. At spot rates (\$2.50/hr per GPU, \$20/hr for the node, 168 hours for a one-week run), the bill is around \$3,360. At on-demand rates and a 3-week schedule for a slightly larger target, it climbs into the low five figures. At a 70B Llama-class campaign on 1,024 H100s for 21 days, the bill is in the millions. The dollar amount scales with the campaign's ambition; the *Spark recipe-lab cost stays at \$1*.

What changes between the small and the large case is the **expected savings**, not the absolute Spark cost. At a 50% wrong-pick rate without prior signal — conservative for blind architectural search — the Spark filter prevents one wrong final-training booking per two campaigns in expectation. For the small case, that is \$1,680 of expected savings per Spark dollar spent. For the medium case, \$7,055. For the 70B Llama-class case, the savings cross seven figures. The ratio of expected savings to Spark cost is ~1670× for the small campaign and ~7000× for the medium one.

```
$ python3 evidence/cost_arithmetic.py | jq .expected_value_argument
{
  "wrong_pick_rate": 0.50,
  "expected_loss_from_blind_booking_usd": 1680.00,
  "spark_recipe_search_cost_usd": 1.01,
  "expected_savings_per_campaign_usd": 1678.99,
  "ratio_savings_to_spark_cost": 1670
}
```

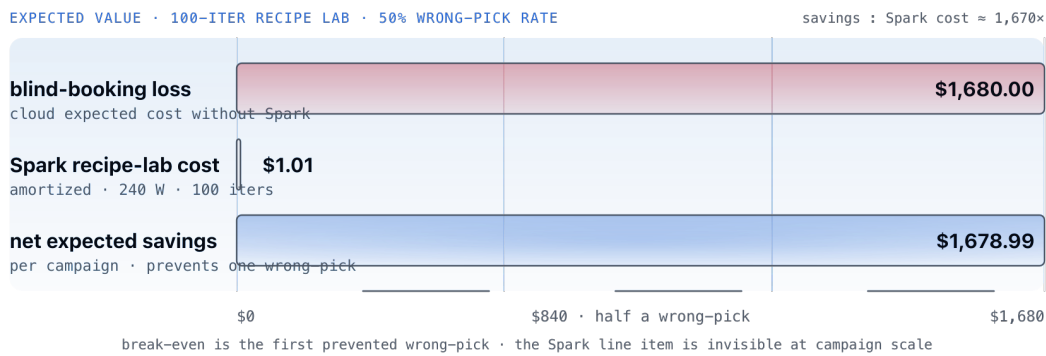


Figure 2.2: Expected-value waterfall for the small campaign. Three horizontal bars at the same scale. Top bar — expected loss from blind cloud booking — extends to \$1,680, rendered in muted red as the cost you absorb without the Spark filter. Middle bar — Spark recipe-lab amortized cost — is a single-pixel sliver at \$1.01, barely visible at the same scale. Bottom bar — expected savings per campaign — extends to \$1,679, the accent, rendered as the value the Spark keeps when one wrong-pick is prevented. The savings bar is nearly identical in length to the loss bar because the Spark cost is so small the difference is invisible at a campaign-budget scale. Right margin annotates the ratio of savings to Spark cost as approximately 1,670 times.

There is no scenario inside this arithmetic where buying a Spark and using it as a recipe lab loses money against the alternative of going to the cloud blind. The break-even is reached after the very first prevented wrong-pick.

(↵) The break-even, on the back of a napkin

Spark amortized cost per 100-iter recipe lab: \$1.01. Expected loss from blind cloud booking at 50% wrong-pick rate × \$3,360 small campaign = \$1,680. Savings per campaign = \$1,679. At one campaign per quarter, the Spark earns back its \$5,000 sticker price in roughly *9 months* of derisking — before any of its other roles (inference rig, fine-tuning sandbox, kernel-envelope test bench) are counted. The break-even point is the *first* prevented wrong-pick.

2.6 Verification — what success looks like on the Spark for a cloud-bound workload

Success in Phase 1 is a 100-row `trajectory.jsonl` file under `articles/autoresearch-agent-loop/evidence/` with `val_bpb` deltas large enough to rank the candidates. Success in Phase 2 is three loss curves drawn over 1,000 steps each with no late-onset divergence. Both artifacts are plain text. Both are reproducible. Both are what you bring to the cloud-vendor billing portal as the justification for booking the node.

What success specifically does *not* look like on the Spark is convergence. You are not training the model on the Spark — you are searching for the architecture you will train on the cloud. The Spark’s job is to give you a defensible reason to skip 90 candidates and pay attention to the other 10. The cloud’s job is to take your three finalists and produce one trained model. Neither machine is being asked to do the other’s job.

Hardware-aware verification: for Phase 1, watch GPU utilization during a representative iteration with `nvidia-smi -lms 500` and confirm the Spark is doing real work (steady ~95% utilization, ~80 GB of unified memory engaged). For Phase 2, log loss every 50 steps and plot the three candidates side-by-side; any candidate that goes flat or diverges before step 800 has just told you something the cloud would have charged \$300 to learn. For Phase 3, the loss curve in the first 24 hours of cloud training should land within ~5% of the Spark stress-test extrapolation. If it doesn’t, your proxy substitution is broken — investigate before letting the cloud node burn through the rest of the week.

2.7 Tradeoffs and gotchas

The argument here has four important caveats. None of them break it; all of them affect how you run a real campaign.

⚠ Spark FP8 numerics aren't bit-exact with H100 FP8

Hopper FP8 (E4M3/E5M2 with Hopper's calibration scheme) is not bit-exact to Blackwell FP8, and Spark's GB10 is not bit-exact to either. Throughput delta is sub-2%; accuracy delta is essentially zero on most architectures, but a recipe whose convergence depends on a specific FP8 calibration choice can silently diverge between Spark and the cloud target. Defense: budget one cloud-side day at the front of Phase 3 for an FP8 sanity check before committing the full week's spend.

FP8 numerics differ between GB10 and H100. Hopper FP8 (E4M3 and E5M2 with their specific calibration scheme) is not bit-exact to Blackwell FP8 — and the Spark's Blackwell GPU is not bit-exact to either Hopper FP8 or to the next-gen B200. In practice the throughput delta is under 2% and the accuracy delta is essentially zero on most architectures, but a recipe that depends on a specific FP8 calibration choice should be validated with one Phase 2 run on the cloud target before committing the full Phase 3 budget. The fix is to budget for an extra cloud-side day at the front of Phase 3 — call it the “FP8 sanity check” — and treat it as a tax against the Spark filter's savings.

📖 Tensor parallelism vs pipeline parallelism

Tensor parallelism (TP) shards a single weight matrix across GPUs within a node — every matmul triggers a collective comm; needs NVLink bandwidth. *Pipeline parallelism* (PP) splits *layers* across GPUs and passes activations forward; cheaper inter-node comms but pipeline bubbles eat efficiency. A typical 7B cloud campaign uses TP=2 within an 8-GPU node and DP=4 across; a 70B campaign adds PP=2 across nodes. The Spark's TP=PP=1 single-GPU mode validates *neither*.

Single-GPU sweep does not exercise multi-GPU parallelism at all. If your final cloud training uses tensor parallelism (TP=2, TP=4) or pipeline parallelism (PP=2, PP=4), the Spark is not validating those code paths. Plan a cloud-side “TP=4 / PP=2 sanity check” for one day before opening the rest of Phase 3 — this is not optional, it is how you avoid discovering on day three of a week-long run that your tensor-parallel attention is wrong.

Memory-bound workloads behave differently with PCIe between GPUs. The Spark has 128 GB of unified memory and zero PCIe between CPU and GPU; a cloud H100 has 80 GB per GPU with PCIe 5 between them. KV cache-heavy or very-long-sequence workloads can hit memory-traffic patterns the Spark cannot reproduce. The defense is the same as for FP8: budget one cloud-side day to confirm that the Spark-measured throughput shape transfers under the cloud's memory topology.

The agent loop's 60-step filter is not exhaustive. Even with Phase 2's 1,000-step stress test, the longest a candidate runs before the cloud commit is roughly 1/100th of its eventual cloud schedule. Some failure modes — gradient instabilities from particular LR-schedule + batch-size interactions, late-onset loss spikes from data ordering — only manifest at scale. The Spark filter cuts the wrong-pick rate; it does not eliminate it. The expected-savings math in the arithmetic table assumes a 50% post-filter wrong-pick rate, which is conservative; in practice with both Phase 1 and Phase 2 running cleanly you can expect closer to 10–20%. That makes the savings math even more lopsided.

🔍 GOING DEEPER

- [Chinchilla scaling laws \(Hoffmann et al., 2022\)](#) — the 20-tokens-per-parameter result that sets Phase 3's Chinchilla-optimal token budget.
- [Kaplan scaling laws \(Kaplan et al., 2020\)](#) — the prior result Chinchilla updated; still the basis for proxy-substitution justification.
- [gpu-sizing-math-for-fine-tuning](#) — sibling Looking-Beyond-Spark piece that walks the same math for fine-tuning instead of pretrain.

- `autoresearch-agent-loop` — the agent harness Phase 1 reuses; the `trajectory.jsonl` format that becomes portable artifact across campaigns.

2.8 What this unlocks

Three concrete things the reader can do this week with what's in this article.

Run a real recipe-lab session for a cloud booking you've been postponing. If you have a planned pretrain on the calendar but haven't decided what architecture to commit to, copy `evidence/recipe_lab_template.py` into the `autoresearch` agent loop's `evidence directory`, point its `ProxyMenu` at your target shape, and let it run overnight. The next morning you have a `trajectory.jsonl` to take to the cloud-vendor billing portal — and a defensible story about why this architecture is the one to fund.

Compute your own break-even number for the Spark as a recipe lab. Re-run `cost_arithmetic.py` with your actual cloud rates, your actual planned campaign size, and your honest estimate of your wrong-pick rate without the filter. The ratio of expected savings to Spark cost is rarely below 100× and often over 1,000×. Once you have that number for your specific situation, the Spark stops being “the small machine on the desk” and starts being “the line item in the AI infrastructure budget that pays for itself in one campaign.”

Treat the trajectory log as a portable artifact across campaigns. A `trajectory.jsonl` from one campaign — even a campaign you never ran on the cloud — is calibration data for the *next* one. The agent loop in `autoresearch` reads its own history when proposing candidates; over multiple campaigns the proposer learns which knobs predict which gains. The first recipe-lab session you run is the most expensive (it has to discover everything from scratch). The fifth is essentially free (the proposer already knows what to try). This is the per-Spark, per-builder version of what the Llama-class scaling-laws papers did at company scale — and it is exactly the kind of thing the `guardrails for code generation` exist to keep honest as the trajectory grows across many sessions.

2.9 State of the apps — as of this article

The **Looking Beyond Spark** thread is now three articles long: [the 100B Nemotron sizing piece](#) (the first in the series), [the layman recap of the autoresearch loop and its 4-tier training roadmap](#), and this one. All three follow the same pattern — work the cloud-side arithmetic on the Spark, present the math as the artifact, and treat the cloud spend as a downstream consequence of an upstream decision. The other three arcs are unchanged: Second Brain has [its four pieces](#), the LLM Wiki arc remains unstarted, and the Autoresearch arc sits at five published pieces ([NeMo Framework](#), [the baseline training loop](#), [the Curator data-prep envelope](#), [the agent loop](#), [the guardrails](#), and [the layman recap](#)) with four to go.

The Spark is now wearing a third hat in this blog. First it was an inference rig. Then it was a training rig. With this article it is also a recipe lab — the small machine on the desk that decides what the big machines in the cloud get to do. That is, it turns out, the most expensive role you can possibly assign to a \$5,000 box, because the cost being measured is not the box's electricity but the cloud booking it gates. The Spark on the desk pays for itself on the *first* prevented wrong-pick. After that, every Spark dollar is making the cloud bill smaller.

📦 Same recipe-lab logic, frontier campaign sizes

The 1,670× savings ratio is the small-campaign case (\$3K cloud commit). Scale up the campaign and the ratio scales with it. A 70B Llama-class run at 1,024× H100 80 GB for 21 days is roughly \$3M; the Spark's \$1 recipe-lab cost gates ~\$1.5M of expected savings at a 50% wrong-pick rate — a 1.5-million-times leverage ratio. A 400B-class frontier run on 4,096× B200s for two months is in the tens of millions; the Spark's role doesn't change, but the savings cross eight figures. The Spark is the one piece of hardware whose ROI *grows* with every other rack you don't yet own.

3. Keep reading

This is a free sample chapter from *AI Research on NVIDIA DGX Spark* by Manav Sehgal.

The full edition has 63 chapters across 9 parts. You can buy the ebook and PDF at leanpub.com/ai-research-nvidia-dgx-spark.