**EPISODE #39**

# SLMs: Precision Power Beyond LLMs

Published December 09, 2025 • Runtime: 22:40

https://myweirdprompts.com/episode/small-langugage-models/

## EPISODE SYNOPSIS

Everyone's heard of Large Language Models, but what about their unsung counterparts? This episode unpacks Small Language Models (SLMs), revealing why they're not just "mini LLMs" but specialized, purpose-built powerhouses. Herman and Corn explain how SLMs are transforming AI workflows, enabling modularity and efficiency, from orchestrating complex tasks as "planning models" to powering AI directly on edge devices, unlocking new realms of privacy and real-time processing. Discover the crucial role these nimble AIs play in a world dominated by giants, proving that sometimes, smaller truly is smarter.

# TRANSCRIPT

## Corn

Welcome, welcome, welcome to My Weird Prompts! I'm Corn, your perpetually curious co-host, and as always, I'm joined by the ever-insightful Herman. How's it going, Herman?

## Herman

All good, Corn. Ready to dive deep into another fascinating prompt from our producer, Daniel Rosehill. He really knows how to unearth the hidden gems of the AI world.

## Corn

He absolutely does! And today's topic, Herman, is one that I think a lot of people overlook, even as they're bombarded with news about AI every single day. We're talking about Small Language Models, or SLMs. Everyone's heard of LLMs, but what exactly are these smaller counterparts, and why should we care?

## Herman

That's a fantastic question to start with, Corn. Because while Large Language Models like GPT-4 or Gemini have captured the public imagination, the unsung heroes often operating behind the scenes, or in very specific niches, are these Small Language Models. And the prompt specifically asks what's out there besides the big players and what roles they play in today's AI workflows. The interesting thing is, they're not just "smaller versions" of LLMs; they often serve entirely different, but equally crucial, functions.

## Corn

Okay, hold on, Herman. I mean, my first thought, and I imagine many listeners' first thought, is that an SLM is just... a smaller LLM, right? Like, a compact version for when you don't need the full powerhouse. A Honda Civic compared to a monster truck. Is that not the gist of it?

**Herman**

Well, I appreciate the analogy, Corn, but I'd push back on that actually. That's a common oversimplification. While some SLMs are indeed distilled or quantized versions of larger models, many are designed from the ground up with specific, constrained tasks in mind. They're not just "less powerful monster trucks"; they're often highly specialized tools, like a precision screwdriver versus a sledgehammer. Each has its job.

**Corn**

Hmm, a precision screwdriver versus a sledgehammer. Okay, I like that distinction. So, it's not just about scale, it's about *purpose* and *design*. Can you give us a concrete example of a genuinely "small" language model, perhaps one that's been around for a while, that illustrates this point? Daniel mentioned BERT in his prompt.

**Herman**

Absolutely. BERT, or Bidirectional Encoder Representations from Transformers, is a prime example. It was introduced by Google in 2018, and it was a game-changer for natural language understanding at the time. While it's dwarfed by today's multi-billion parameter LLMs, BERT is still widely used in production for tasks like sentiment analysis, text classification, and named entity recognition. Its design is optimized for understanding context within a sentence, rather than generating long, coherent prose. It's incredibly efficient for those specific tasks.

**Corn**

So, BERT isn't trying to write a novel; it's trying to figure out if someone's tweet is positive or negative, or if a word in a sentence is a person's name.

**Herman**

Precisely. And that's a key differentiator. A truly small language model is often built with a smaller architecture, fewer parameters, and trained on a more focused dataset for a particular domain or task. This makes them faster, less resource-intensive, and often more accurate for *that specific task* than an LLM trying to do everything.

**Corn**

Okay, but then what about the "quantization" aspect you mentioned? Because I've definitely seen terms like "quantized models" or "quantization of LLMs" flying around. How does that fit into the SLM picture, or does it?

**Herman**

That's where it gets a bit nuanced, and it's important to make the distinction. Quantization is a technique used to reduce the size and computational cost of an *existing* model, usually a larger one, by representing its parameters with fewer bits of information. Think of it like compressing a high-resolution image into a lower-resolution one. You still have the same underlying image, just with less detail and a smaller file size.

**Corn**

So, a quantized LLM is still fundamentally an LLM, just on a diet?

**Herman**

Exactly. A quantized LLM is an LLM that has undergone a process to make it smaller and faster, often suitable for deployment on edge devices or with less powerful hardware. It's still trying to do the same things as its larger parent model, albeit with some potential loss of fidelity. A true SLM, on the other hand, might have been designed with a modest parameter count from the very beginning, optimized for a specific job without necessarily starting as a massive model.

**Corn**

That makes a lot more sense. So, a quantized Mistral 7B is an LLM that's been shrunk down, but a purpose-built model like BERT is a genuinely small language model by design. And you mentioned Hugging Face, which Daniel did too. It sounds like a real treasure trove for these kinds of models.

**Herman**

It absolutely is. Hugging Face is essentially the GitHub for AI models. It's a massive repository where researchers and developers share models, datasets, and even entire development environments. If you're looking for open-source AI models, from the colossal to the miniscule, that's often the first place to check. It's democratizing access to AI, which is a fantastic thing. You'll find thousands of fine-tuned BERT variants, specialized text classifiers, code models, and so much more that are far from the headline-grabbing LLMs.

**Corn**

It's a bit of a model jungle, as Daniel put it in his prompt, but in a good way. Like a vibrant ecosystem where everything has its place. So, moving beyond just BERT, what other categories or examples of these truly small, purpose-built models are gaining traction? You talked about "planning models" in the prompt.

**Herman**

Right. Beyond the well-known ones like BERT or specialized sequence-to-sequence models for tasks like machine translation, we're seeing an emergence of what I like to call "accessory models" or "helper models" within larger AI workflows. These are often the "planning models" Daniel refers to. Think of them as specialized internal components in a complex machine.

**Corn**

Like, miniature AI assistants to the main AI?

**Herman**

You could say that. For instance, in a complex RAG—Retrieval-Augmented Generation—system, you might have a small, highly efficient SLM whose sole job is to re-rank search results before they're fed to the main LLM. Or another SLM designed specifically to classify the *intent* of a user's query, directing it to the right tool or sub-system, rather than having the LLM figure that out itself. This is often more reliable and faster.

**Corn**

That's fascinating. So, instead of making the LLM do *everything* – understand the query, search the database, summarize, generate the response – you're breaking it down into smaller, more manageable steps, each handled by a specialized, efficient SLM. It sounds like a modular approach.

**Herman**

Exactly! It's about modularity and distributed intelligence. For example, a "planning model" might receive an initial complex user request like, "Plan a five-day trip to Rome that includes historical sites, excellent food, and a day trip to Pompeii." Instead of the main LLM directly generating the whole itinerary, the planning SLM could first break this down into sub-tasks: `[1. Research Rome historical sites, 2. Find highly-rated Roman restaurants, 3. Plan transportation to Pompeii, 4. Combine into itinerary]`. It then orchestrates which specialized tools or even which *other* SLMs, or eventually the LLM, should handle each sub-task.

**Corn**

Let's take a quick break from our sponsors. Larry: Are you tired of feeling like your brain is just... *there*? Introducing "Cerebral Surge," the revolutionary new brain enhancer that unlocks your mind's latent potential! Our proprietary blend of "bio-luminal frequencies" and "cognitive activators" is scientifically engineered to make you think faster, harder, and... well, just *more*. Side effects may include sudden urges to reorganize your pantry, an inexplicable affinity for abstract art, and occasionally, feeling like you understand quantum physics. Cerebral Surge: Because mediocrity is *so* last century. Buy now!

**Herman**

...Alright, thanks Larry. Anyway, where were we? Ah yes, modularity. Corn, you brought up a good point about orchestration. This is where SLMs really shine, particularly in enterprise or complex AI applications.

**Corn**

So, we're essentially building a team of specialists, where the LLM is the brilliant but somewhat generalist CEO, and the SLMs are the highly efficient, task-specific department heads.

**Herman**

That's a great analogy, Corn. It's about optimizing the entire workflow. LLMs are powerful, but they're also resource-intensive and can be slow. By offloading specific, well-defined tasks to SLMs, you gain several advantages: **speed**, because a smaller model processes data much faster; **cost-efficiency**, as running SLMs is significantly cheaper than constantly querying an LLM; and **reliability**, because a model trained narrowly on a specific task often performs that task with higher accuracy and fewer hallucinations than a general-purpose LLM.

**Corn**

But wait, Herman, aren't we just introducing more complexity by having all these different models talking to each other? More points of failure, more things to manage. For a lot of businesses, simplicity is key, right? They just want one AI to do the job.

**Herman**

I'd push back on that actually. While it adds a layer of architectural complexity, the *operational* benefits often outweigh it. Think about software development. You don't build an entire application as one monolithic block anymore; you break it down into microservices. Each microservice does one thing well, and they communicate via APIs. If one microservice fails, the whole system doesn't necessarily crash, and you can update or scale individual components. SLMs are the microservices of the AI world. They allow for more robust, scalable, and maintainable AI systems.

**Corn**

Okay, I see your point. It's about building resilient systems. So, beyond orchestrating larger workflows, what about their use cases in environments where a huge LLM just isn't practical? Like on your phone, or an IoT device?

**Herman**

Excellent point! That's another critical area: **edge computing**. Because SLMs are so much smaller and less demanding, they can run directly on devices like smartphones, smart speakers, even drones or industrial sensors. This means real-time processing without needing to send data to the cloud, which has huge implications for **privacy**—as sensitive data stays on the device—and for **latency** in applications where speed is paramount. Imagine a smart camera on a factory floor that uses a small vision language model to detect anomalies in real-time, without having to send every frame to a distant server.

**Corn**

So, they're not just accessories to LLMs; they're also enabling entirely new classes of on-device AI applications. That's a huge potential market. I mean, my phone struggles enough with running regular apps, let alone a giant LLM!

**Herman**

Exactly. And the ability to run these models locally opens up opportunities for personalized AI experiences, as the model can be fine-tuned with your specific data without ever leaving your device. We're also seeing SLMs being used for sophisticated content moderation, data governance, and even generating synthetic data for training larger models. The potential is vast.

**Corn**

Alright, we've got a caller on the line. And we've got Jim on the line – hey Jim, what's on your mind? Jim: Yeah, this is Jim from Ohio. And I've been listening to you two go on about these "small models" and "big models" and "quantized" whatever, and frankly, I think you're making a mountain out of a molehill. It just sounds like you're trying to sell me more things. My neighbor Gary, he's always trying to sell me some new gadget for my lawnmower, says it'll make it run better. It never does. And now you're telling me I need a whole team of little AIs just to make one big AI work? Sounds like a lot of extra steps for nothing.

**Herman**

I appreciate your skepticism, Jim, and it's a valid concern about complexity. But the truth is, these small models aren't about selling you *more*; they're about making the overall system *more efficient and effective*. Think of it like a specialized pit crew at a race. Each member has a specific, small task they do incredibly well and quickly, allowing the race car—your "big AI"—to perform at its peak without being bogged down.

**Corn**

Yeah, Jim, it's not about making things complicated for the sake of it. It's about specialization. You wouldn't use a sledgehammer to hammer in a thumbtack, right? These small models are the digital thumbtack hammers. They do their one job precisely and quickly, which actually *reduces* the overall cost and time compared to having a general-purpose tool try to do everything. Jim: Ehh, I don't buy it. Sounds like more processing, more electricity. In my day, if you wanted something done, you just did it yourself. No fancy AI teams. And the weather here in Ohio has been all over the place, one day it's sunny, the next it's raining cats and dogs. Can't trust anything these days. But seriously, this just sounds like over-engineering to me. Why can't the big one just figure it out?

**Herman**

Well, the big one *can* figure it out, Jim, but at a higher cost in terms of computational resources and time. By distributing the workload to specialized SLMs, we're not just doing the same thing with more steps; we're often doing it *better*, *faster*, and *cheaper* for specific parts of the process. It's an optimization strategy.

**Corn**

It really is, Jim. It's like having a dedicated librarian who knows exactly where to find every book, versus asking a general encyclopedia to search its entire contents every time you have a question. The librarian is faster for that specific task. Jim: Still sounds like too much fuss. Thanks for nothing. And my cat Whiskers is giving me that look like she knows something I don't. I gotta go.

**Corn**

Thanks for calling in, Jim! Always a pleasure.

**Herman**

Always… insightful. Well, Jim brings up a common sentiment, Corn. People often crave simplicity. But sometimes, true simplicity and efficiency in complex systems are achieved through elegant modularity.

## Corn

It's a good point, and one worth addressing. So, for our listeners who are maybe working in AI development, or leading teams that use AI, what are the practical takeaways here? How can they actually leverage this understanding of SLMs in their own workflows?

## Herman

For developers and architects, the key takeaway is to start thinking of your AI systems not as monolithic LLM deployments, but as integrated ecosystems. Identify tasks within your workflow that are well-defined and don't require the full generative power of an LLM. Could an SLM handle data cleaning, intent classification, sentiment analysis, or initial data retrieval more efficiently?

## Corn

So, it's about breaking down the problem into smaller, bite-sized pieces and then matching the right tool – big or small – to the right job. That makes a lot of sense from an engineering perspective.

## Herman

Exactly. This leads to several benefits: **reduced operational costs** because you're not paying for expensive LLM inference on every single request; **improved latency** for user-facing applications; **enhanced privacy** if certain SLMs can process sensitive data locally; and **greater robustness** because if one small model fails, the entire system isn't necessarily crippled. You also get **easier fine-tuning** since a smaller model trained on a narrower dataset is often more amenable to specific adaptations without "catastrophic forgetting."

## Corn

But isn't there a risk of getting lost in the "model jungle" then, as Daniel mentioned? With so many models out there, how do you even choose the right one for a specific task?

**Herman**

That's a fair challenge. It requires a good understanding of your specific problem domain and familiarity with resources like Hugging Face, where you can explore and benchmark various models. It also means potentially investing in more sophisticated orchestration frameworks to manage the interactions between different models. It's a shift in architectural mindset, from "one model does it all" to "a suite of models collaborates to solve the problem." But for certain applications, the payoff in efficiency and performance is significant.

**Corn**

So, for business leaders, it's about recognizing that not every AI problem needs a multi-billion dollar supercomputer. There are often more economical, faster, and sometimes even more accurate solutions in the realm of these specialized SLMs.

**Herman**

Precisely. And for everyone, it's about understanding that the AI landscape is far richer and more diverse than just the headlines suggest. The future of AI isn't just about bigger models; it's also about smarter, more distributed, and more specialized applications of AI, where SLMs play an absolutely critical role.

**Corn**

That's a truly insightful perspective, Herman. It really changes how you look at the whole AI ecosystem. It's not a single towering skyscraper; it's a vast, interconnected city with all sorts of buildings, big and small, each serving a vital purpose.

**Herman**

Well put, Corn. The interplay between LLMs and SLMs, and the continuing innovation in specialized model design, promises an exciting future for AI.

**Corn**

Absolutely. It makes you wonder what other hidden gems Daniel will unearth for us in future prompts. So much to explore!

**Herman**

Always.

**Corn**

And that brings us to the end of another thought-provoking episode of My Weird Prompts. A huge thank you to Daniel Rosehill for sending in such a fascinating prompt and getting us to really dig into the world of Small Language Models.

**Herman**

Indeed. It's a topic I think we'll be hearing a lot more about as AI systems mature.

**Corn**

For sure. You can find "My Weird Prompts" on Spotify and wherever else you get your podcasts. Make sure to subscribe so you don't miss an episode. We love hearing from you, so if you have any weird prompts of your own, send them our way! Until next time, I'm Corn.

**Herman**

And I'm Herman.

**Corn**

Stay curious!