

## MY WEIRD PROMPTS

Podcast Transcript

### EPISODE #29

# The Multimodal Audio Revolution: A Screen-Free Future?

Published December 07, 2025 • Runtime: 25:47

<https://myweirdprompts.com/episode/audio-multimodal-vs-stt/>

## EPISODE SYNOPSIS

Welcome to "My Weird Prompts"! This episode, Corn and Herman dive into producer Daniel Rosehill's fascinating concept of "audio multimodal modality," which he champions as the next major wave of speech technology. Is this advanced AI, capable of understanding context, tone, and performing complex tasks from simple audio prompts, truly set to displace traditional speech-to-text models entirely? Herman unpacks how these multimodal systems go beyond mere transcription to offer a profound shift towards screen-free work, enhanced accessibility, and intelligent content creation. However, he also challenges Daniel's bold prediction, exploring where classic STT will continue to play a vital, specialized role due to factors like cost, data integrity, and real-time demands. Join them as they explore the potential and practicalities of this groundbreaking evolution in audio AI, asking if we're on the cusp of a truly screen-free future, or if specialized tools will always have their place.

# TRANSCRIPT

## Corn

Welcome back to "My Weird Prompts," the podcast where we dive deep into the fascinating questions and ideas sent to us by our very own producer and creator, Daniel Rosehill. I'm Corn, and I am absolutely thrilled to be here as always with my brilliant co-host, Herman.

## Herman

And I'm Herman. Today, Corn, Daniel has sent us a prompt that's not just about a nascent technology, but about a potential paradigm shift in how we interact with information and, crucially, how we might reduce our reliance on traditional screen-based work. It's a big one.

## Corn

Oh, I'm already hooked! Daniel mentioned he's particularly excited about this area. He describes it as the "audio multimodal modality," which is a bit of a mouthful, but he believes it's not getting the attention it deserves compared to, say, video or images in the AI space. He thinks this could be the next major wave of speech technology, potentially displacing classic speech-to-text models entirely. Herman, is Daniel onto something that big?

## Herman

Absolutely, Corn. Daniel is not just onto something, he's highlighting a fundamental evolution in how AI processes and understands audio. What he refers to as "multimodal audio modality" or more simply, multimodal audio models, represents a significant leap beyond traditional automatic speech recognition, or ASR, which is essentially what classic speech-to-text, or STT, models do.

## Corn

So, let's unpack that a bit. What exactly \*is\* a multimodal audio model, and how is it different from just, you know, transcribing words? Because I use transcription services all the time, and they just give me the text. Is this a fancy new name for the same thing?

## Herman

Not at all, Corn. Think of classic STT as a highly specialized but narrow-focused tool. Its primary function is to convert spoken words into written text. It excels at recognizing phonemes, assembling them into words, and then sentences. However, it typically lacks an inherent understanding of the \*meaning\* or \*context\* of those words beyond linguistic patterns. It's like having a brilliant stenographer who can capture every word perfectly but doesn't necessarily grasp the nuances of the conversation or the speaker's intent.

## Corn

Okay, so it's good at \*what\* was said, but not \*why\* or \*what it means\*?

## Herman

Precisely. Now, enter multimodal audio models. These are engineered to integrate and process information from multiple types of data, or "modalities," simultaneously. In the context of audio, this means they don't \*just\* transcribe the speech; they also analyze the audio for other cues, such as tone, emotion, pauses, speaker identification, and importantly, they couple that with advanced natural language understanding, or NLU, and even natural language generation, NLG. The "multimodal" aspect here specifically refers to combining the raw audio signal with linguistic processing and a broader understanding of context, often through a text prompt.

## Corn

So, it's like a super-smart assistant who listens to what you say, \*and\* how you say it, \*and\* what you mean by it, and then can act on that deeper understanding? That's wild. Daniel specifically mentioned his experience with a tool like Whisper for transcription, and then having to use other tools or manual processes to clean up and format the text. He gave an example of trying Gemini on an experimental one-hour transcription job, asking it not just to transcribe but also to "clean it up a little bit" and "add diarization." He saw this as a huge leap.

## Herman

That's a perfect example, Corn. In a classic ASR or STT pipeline, if you wanted diarization—identifying who said what—or if you wanted to clean up filler words, correctly punctuate, or summarize, those would typically be separate, downstream processes. You might need different models, or you'd apply a series of rules or even manual editing. But with multimodal models, particularly those leveraging large language models, or LLMs, you can embed those instructions directly into your natural language prompt alongside the audio.

## Corn

So, instead of "transcribe this audio" and then a separate command of "now, identify the speakers" and "now, remove the ums and ahs," you just say "transcribe this audio, clean it up, and tell me who said what"? That sounds incredibly efficient.

## Herman

It is. Daniel's observation about Gemini is spot-on. The ability to integrate complex processing instructions directly into the audio prompt, in natural language, is a game-changer. It leverages the contextual understanding and generation capabilities of the LLM to interpret not just the spoken words, but the *\*intent\** behind the prompt for how that audio should be processed and presented. This moves beyond simple transcription to intelligent summarization, content restructuring, and even creative output, all from a single input stream.

## Corn

And Daniel specifically highlighted that this combination of accurate transcription—which is where STT models excel—with the ability to then process that raw, often messy, spoken-word text, is where the real value lies. He mentioned how raw transcription often misses punctuation, paragraph breaks, and includes all those natural speech quirks like "um" and "uh." This new approach can fix all that.

## Herman

Indeed. The text output from raw STT is inherently different from typed text. It's often stream-of-consciousness, laden with disfluencies, false starts, and lacks the structural elements we expect in written communication. A multimodal model, by understanding the *\*intention\** to create a polished output, can intelligently infer punctuation, structure, and even refine word choices to produce a more coherent and readable text. It's not just a transcription; it's a transformation.

## Corn

This is where Daniel sees a totally novel form of sharing information digitally, one that could significantly reduce the need for people to be tethered to their desks. He argued that audio is much more versatile, allowing you to create content from anywhere, as long as you can speak into a microphone. This sounds like a productivity play on one hand, but he also hinted at something more powerful and wonderful.

## Herman

That dual perspective is key, Corn. On the productivity side, imagine knowledge workers, educators, or creatives who spend hours typing emails, reports, or blog posts. With these advanced multimodal audio models, they could simply speak their thoughts, insights, or instructions, and the AI could not only transcribe it accurately but also format it, summarize it, extract key action points, or even draft a coherent document, all based on a high-level natural language instruction. This dramatically reduces the friction between thought and documented output.

## Corn

So I could dictate a complex email, say "Herman, draft an email to the team summarizing our last podcast episode, highlight the key points about multimodal audio, and include a call to action for listeners to check us out on Spotify, but make it sound friendly and professional," and it would just \*do\* it?

## Herman

In principle, yes. And that's not just a productivity gain; it's also about accessibility and empowerment. Daniel mentioned this, too. For individuals who find typing difficult, whether due to physical limitations, dyslexia, or simply preference, this technology could unlock new avenues for communication and creation. It lowers the barrier to entry for content production significantly. This is the "more powerful and wonderful" aspect Daniel was alluding to. It's about empowering people to leverage their most natural form of communication—speech—to interact with and shape the digital world.

## Corn

That's a huge point about accessibility. It's not just about convenience for those who \*can\* type, but enabling those who \*can't\* to participate more fully. Daniel's prediction, then, is that this capability is so strong that these multimodal audio models will entirely \*displace\* classic speech-to-text. He asked us to challenge him on that. So Herman, is there really no longer a need for basic, raw STT? Will it truly become obsolete?

## Herman

That's the crux of Daniel's challenge, and it's a very valid question. While the capabilities of multimodal audio models are undeniably impressive and represent the future for many applications, I would argue that classic STT, or at least its underlying function, will likely persist, though perhaps in a more specialized or foundational role. It's not necessarily displacement as much as re-contextualization.

## Corn

Okay, so you're pushing back a bit. Where do you see classic STT still holding its own against these newer, more powerful models? Give me some concrete use cases or reasons.

## Herman

Let's consider several factors, Corn. First, **cost and computational resources**. Multimodal models, especially those integrating large language models, are significantly more resource-intensive to run. They require powerful GPUs and substantial processing power. For simple, raw transcription where only the text stream is needed, a lightweight, classic STT model can be much faster and cheaper to operate, especially at scale. Think of high-volume transcription services where only the raw words are needed for archival or basic search, without immediate interpretation or refinement.

## Corn

So if I just need a transcript of an hour-long meeting for my records, and I'll clean it up myself later, or I'm happy with the rough output, the simpler, cheaper option might still win out?

## Herman

Precisely. And this ties into the next point: **Specialized accuracy and raw data integrity**. In certain domains, like legal proceedings, medical dictation, or financial compliance, the absolute fidelity of the raw spoken word is paramount. Any "cleaning up" or "interpretation" by an AI, no matter how sophisticated, could introduce unwanted bias or alter the original meaning. In these scenarios, a raw, unedited, highly accurate STT output, perhaps with timestamps, serves as the definitive record. The human expert then performs the interpretation or refinement, often based on specific domain knowledge that an AI might struggle to fully grasp or apply without hallucination.

## Corn

That's a good point. The AI cleaning up my podcast summary is one thing, but if it's medical notes, I definitely want the doctor's exact words, even if they're a bit messy.

### Herman

Exactly. Then there's **privacy and data security**. For highly sensitive audio data, organizations or individuals might prefer to process it using local, on-device STT models that never send the audio data to a cloud-based multimodal service. While multimodal models can also be run locally, they typically require more sophisticated hardware, making a simple, local STT model a more accessible and secure option for basic transcription needs without cloud reliance.

### Corn

So if I'm recording something extremely confidential, I might want a simpler model that doesn't send my data off to a big server farm for complex processing.

### Herman

That's a strong consideration for many. Another factor is **latency and real-time applications**. For real-time applications where immediate text output is critical, such as live captioning or voice command interfaces that don't require complex semantic understanding, classic STT models can offer lower latency. They are optimized for speed in converting speech to text, whereas multimodal models might introduce slightly more delay due to the additional layers of processing and understanding.

### Corn

Ah, so for live captions on a broadcast, or talking to my smart home device, I just need it to get the words right *now*, not necessarily understand the deeper meaning of my request to turn on the lights.

### Herman

Right. The simple, direct action. Also, we must consider **system complexity and integration**. For developers building applications, sometimes integrating a simple STT API that just returns text is far easier and lighter-weight than integrating a full multimodal model that has many more parameters and potential outputs. It's about choosing the right tool for the job. If the job is just "get me the words," classic STT remains straightforward.

## Corn

So, it's not a complete "either/or" situation, then. More of a spectrum of tools for different purposes. Herman, what about the potential for multimodal audio models to introduce new types of errors or biases compared to classic STT? If they're interpreting and cleaning up, could they misinterpret or make assumptions that a raw transcription wouldn't?

## Herman

That's a crucial point, Corn, and a significant challenge for multimodal models. While classic STT can make transcription errors—mishearing a word, for instance—its errors are typically phonetic or lexical. The output is usually a faithful, albeit sometimes imperfect, representation of the \*literal\* spoken words. Multimodal models, however, introduce a layer of interpretation. When you ask it to "clean up" or "summarize" or "format," it's applying its learned understanding of language and context, which can introduce: 1. **Hallucinations**: The model might "invent" information or rephrase something in a way that subtly changes the original meaning, even if it sounds coherent. 2. **Bias**: If the training data contains biases in how certain speech patterns, accents, or topics are interpreted, the "cleaned up" or summarized output could inadvertently perpetuate those biases. 3. **Loss of Nuance**: In attempting to streamline or summarize, subtle inflections, implied meanings, or specific phrasing that was important to the original speaker might be lost. 4. **Over-correction**: The model might be overly aggressive in removing "ums" and "ahs" or restructuring sentences, inadvertently losing the natural cadence or intent of the speaker.

## Corn

So, in an effort to make it sound "better," it might actually distort the original message or add things that weren't there. That's a real risk, especially if you're not carefully reviewing the output.

## Herman

Exactly. This highlights the ongoing importance of the "human in the loop." Even with advanced multimodal models, human review and editing remain critical for high-stakes or high-fidelity applications. The models can do 90% of the heavy lifting, but that last 10% often requires human judgment, domain expertise, and an understanding of the specific communication goal.

## Corn

So Daniel's initial excitement about the "raw, very long text that arrives from a transcription" being processed and cleaned up seamlessly is valid, but we need to be aware of what might be lost or changed in that process. It's not just a trivial "use case," as he put it, to simplify that process. It's a powerful one, but with caveats.

## Herman

It's powerful precisely *because* it tackles that common pain point, but the "simplifying" aspect doesn't necessarily mean "perfecting" or "preserving all intent." It's about automating a step that previously required significant human effort, allowing humans to focus on higher-order tasks like factual verification or strategic refinement.

## Corn

This makes me think about the broader implications. If Daniel is right, and this changes how we share information digitally, what does that mean for how we even *think* about creating content? We've been so text-centric for so long.

## Herman

That's the truly transformative potential, Corn. We're moving towards an era of ambient computing and more intuitive human-computer interaction. Imagine a world where your thoughts, spoken naturally, can seamlessly translate into actionable insights, beautifully formatted documents, or even multimedia content. This isn't just about dictation; it's about a conversational interface to your digital life. For example, a marketing professional could brainstorm ideas aloud during their commute, and the multimodal model could automatically draft social media posts, blog outlines, or even initial campaign strategies from that spoken input. An architect could walk through a construction site, describing observations and instructions, which are then automatically transcribed, diarized, categorized by project, and assigned as tasks to team members.

## Corn

That's not just productivity; that's fundamental shift in how work gets done and how we share knowledge. It decouples the act of creation from the act of typing or formatting.

### Herman

And it extends to consumption as well. Imagine podcasts or lengthy audio lectures not just being transcribed, but dynamically summarized, key takeaways extracted, or even translated into different languages or formats on demand, all through multimodal AI. This could democratize access to information even further. Daniel mentioned the reduction of time spent at desks because audio is more versatile. That's a huge potential societal impact.

### Corn

So, while classic STT might still have its niche for raw, uninterpreted data, the future of general content creation and knowledge management, especially for those looking to move beyond the keyboard, is definitely in these multimodal models. Daniel's prediction isn't about STT disappearing, but its primary \*use case\* shifting dramatically.

### Herman

Precisely. It will become the foundational layer that multimodal models often leverage, much like a raw image sensor captures data that is then processed and enhanced by a sophisticated camera's software. The raw data is still there, still valuable, but the consumer-facing output will increasingly come from these more intelligent, interpretive systems.

### Corn

Herman, this has been an incredibly insightful discussion. It sounds like Daniel's excitement is totally warranted, and his prediction about the "next wave" of speech tech is likely accurate, though perhaps "re-contextualization" rather than outright "displacement" for classic STT.

### Herman

I would agree with that nuance. Multimodal audio models represent a profound evolution, but the underlying utility of accurate, raw speech-to-text remains, albeit for specific, often foundational, purposes. It's an exciting time to be observing this space.

**Corn**

I mean, the ability to simply talk and have an AI process, understand, and then act on those instructions, even cleaning up your ramblings, it really is powerful. Daniel, thank you for sending us such a thought-provoking prompt about the future of speech technology. It's clear why you're so interested in this evolving space.

**Herman**

Indeed, Daniel. Your insights into the practical applications and the deeper implications are always appreciated. You've given us plenty to chew on here.

**Corn**

And to our listeners, if you're as fascinated by these weird prompts and their deep dives as we are, make sure to find "My Weird Prompts" wherever you get your podcasts, including Spotify. We love exploring these cutting-edge topics.

**Herman**

We do. And we look forward to next time.

**Corn**

Until then, stay curious!

**Herman**

Goodbye.