**◎ ChatGPT**

# Azure Custom Voice Services and Capabilities

Azure offers powerful speech AI services that enable creation of custom synthetic voices from user-provided recordings. This report provides a detailed overview of these services – starting with **Custom Neural Voice** – and related text-to-speech (TTS) and speech-to-text (STT) capabilities. It covers available services, how custom voice creation works, pricing details, prerequisites (including responsible AI requirements), and real-world case studies (virtual assistants, accessibility, dubbing, etc.) demonstrating Azure's custom voice features.

## Azure Custom Neural Voice (CNV)

**Custom Neural Voice** (CNV) is an Azure AI service feature that lets you create a unique, natural-sounding synthetic voice by training on recordings of a target voice [1]. It is part of Azure Cognitive Services' Speech service and uses advanced neural text-to-speech technology. CNV is ideal for developing a branded voice or character voice that isn't available among Azure's pre-built voices [2]. By providing your own audio samples and transcripts as training data, you can generate a one-of-a-kind voice font that closely matches the original speaker's tone, mannerisms, and speaking style [1] [3].

**How it works:** To create a custom voice, you use the **Azure Speech Studio** portal. The process involves collecting high-quality voice recordings and their transcripts, uploading this data to a Speech Studio project, and then initiating training of a custom model [3]. Azure's neural TTS engine uses these samples to learn the voice's characteristics (phonetics, prosody, accent, etc.), producing a synthetic voice model. Once trained, the custom voice is deployed to an endpoint for usage in TTS – allowing your applications to convert arbitrary text into speech in that custom voice [3]. Microsoft recommends careful planning of the voice persona and script coverage in the recordings to ensure the model captures the desired style and vocabulary [4]. Creating a high-quality custom voice may require dozens of minutes of recorded speech; Azure supports **Custom Neural Voice Professional** projects for highest quality (using studio-recorded audio) and also offers a lighter option for evaluation called **Custom Neural Voice Lite** [5].

**"Personal Voice" feature:** In 2024, Microsoft introduced *Personal Voice* as a streamlined approach to custom voices for individuals [6]. Personal Voice allows users to create an AI clone of their own voice with a much smaller recording sample (as little as a 1-minute sample for an initial model) [6]. This feature is built on the CNV technology and is intended for personal or limited use cases – such as users wanting their voice in a virtual assistant or for accessibility purposes. Like CNV, the personal voice feature runs on Azure's speech service; users can record a short script, and Azure will produce a custom voice model that sounds like them [6] [7]. This model can then be integrated into apps (e.g. a voice assistant or translator) so that the app speaks in the user's own voice. Personal Voice, while easier to set up, still abides by the same responsible AI guardrails (limited access and required consent – see below).

## Azure Text-to-Speech (Speech Synthesis)

Azure's **Text-to-Speech (TTS)** service converts text into spoken audio using a variety of synthetic voices. Out-of-the-box, Azure provides a large selection of pre-built voices (including **Neural** voices in many languages that sound highly natural) [2] . Developers can call the TTS API or SDK to synthesize speech from text using these standard voices, which are tuned for general use. The neural TTS engine supports rich SSML (Speech Synthesis Markup Language) tags to control pronunciation, speaking rate, pitch, pauses, and even **speaking styles** (like cheerfulness, empathy, newscast style, etc., for supported voices). Azure continually expands its voice library, including **emotional tone** and **style** controls, to make synthesized speech more expressive and human-like.

The TTS service's **voice synthesis capabilities** also extend to custom voices. A custom voice created via CNV becomes accessible through the same TTS API, but as a private voice endpoint. This means your applications can call Azure's TTS using your custom voice, enabling the app to "speak" with that unique voice. The underlying technology for both standard neural voices and custom voices is the neural speech synthesis model, which produces highly realistic prosody and pronunciation [8] . In fact, CNV uses a multilingual, multi-speaker universal model as a base, which can be fine-tuned to the target voice; this approach even allows cross-language voice adaptation (e.g., using an English voice to speak other languages) and a wide range of speaking styles [9] .

Azure's TTS supports **nearly 140 languages and dialects** (as of recent updates) and over 400 voices, covering different genders, ages, and styles. If a truly unique voice is not required, these pre-built voices often suffice for applications; they can be a faster way to integrate TTS (with no training needed) and are available instantly for use [2] . However, for branding or personalization (e.g., a virtual assistant with a signature voice, a distinctive character voice, or mimicking a specific person's voice with permission), Custom Neural Voice provides the path to build that tailored voice.

**Voice output and quality:** Azure's neural voices (including custom ones) are known for their natural intonation and fluidity, using deep learning to handle prosody (tone, stress, timing) in a way that avoids the choppy or robotic sound of older-generation TTS [8] . The service can handle real-time streaming synthesis (for interactive scenarios like voice agents) as well as **batch synthesis** (converting large text documents or books to audio files). A special *Long Audio* API is available for high-quality offline generation of lengthy content (e.g. audiobooks or narration), using more extensive processing for optimum quality – this incurs different pricing (discussed below).

## Azure Speech-to-Text (Speech Recognition)

Azure's **Speech-to-Text (STT)** service (also called Speech Recognition) transcribes spoken audio into text. It supports real-time transcription (streaming microphone input and getting text in milliseconds) as well as batch transcription of audio files. The service is capable of recognizing speech in many languages and can handle various scenarios (conversational input, dictation, etc.). In addition to raw transcription, Azure STT can provide features like automatic punctuation, profanity filtering, speaker diarization (distinguishing between different speakers in multi-speaker audio), and **language identification** (detecting the spoken language) [10] .

Similar to TTS, the speech recognition service offers **pre-trained models** for each supported language that work out-of-the-box. These models are trained on large datasets and perform well for general vocabulary and acoustics. However, Azure also enables **customization of speech recognition** to improve accuracy for specific jargon, noise conditions, or speaking styles. Users can create a **Custom Speech** model by providing their own audio data and transcripts to adapt the base model [11] . This customization can involve tuning an **acoustic model** (to better handle specific background noise or accents) and/or a **language model** (to better recognize domain-specific terminology or phrases) [12] [13] . For example, a custom model could be trained to improve recognition of medical terms for a healthcare app, or to better transcribe speakers with a particular regional accent. These custom STT models are then deployed to an endpoint and used in place of the standard model for transcription.

Using the STT service is straightforward via the Azure Speech SDK or REST API: you send audio and receive text results. The **Azure Speech Studio** provides a graphical interface to test speech recognition and evaluate accuracy, as well as to manage custom models (via the Custom Speech portal). In summary, Azure STT covers scenarios from basic voice transcription to advanced, tuned models for high-accuracy recognition in specialized applications [11] .

## Pricing Details for Custom Voice, TTS, and STT

Azure's speech services are billed primarily on a **consumption model** (pay-as-you-go), with costs based on usage. Pricing can vary by region and currency, but here we present representative costs in US dollars (USD). Table 1 summarizes key pricing metrics for speech services (note that Azure also offers discounted **commitment tiers** for large volumes and a free tier for limited usage):

**Table 1 – Azure Speech Services Pricing (Pay-as-You-Go)**

| Service | Pricing (approx. US rates) |
| --- | --- |
| **Speech-to-Text (STT)** – Standard model | **$1.00 per audio hour** transcribed [14] . Billed per second of audio. |
| **Speech-to-Text (STT)** – Custom model | **$1.40 per audio hour** [14] , billed per second, when using a custom speech model. <br>(Additional ~$0.50 per hour per custom model for endpoint hosting may apply) [14] . |
| **Text-to-Speech (TTS)** – Neural voices (standard/prebuilt) | **$16 per 1 million characters** for real-time or batch synthesis [15] . Billed per character of input text (one character is 1 byte, with SSML tags not counted towards quota). |
| **Text-to-Speech (TTS)** – Long Audio (batch synthesis) | **$100 per 1 million characters** for long-form batch synthesis [15] . This higher rate applies when using the Long Audio API for high-quality audio generation (e.g. large documents, audiobooks). |

| Service | Pricing (approx. US rates) |
|---|---|
| **Custom Neural Voice** – Training a custom voice model | **$52 per compute hour** of training time [16] , up to a maximum of 96 compute hours per training (i.e. capped at ~$4,992 for a single training job) [17] . Training time depends on data size and model type (a high-quality professional voice might use 20–40 compute hours for a single-style voice, and up to ~90 for a multi-style voice) [18] . |
| **Custom Neural Voice** – Voice synthesis (usage) | **$24 per 1 million characters** for using the custom voice in speech synthesis (real-time or batch) [16] . This is the runtime cost to generate speech with your custom voice, billed per character of input text (higher than standard voices due to custom model usage). |
| **Custom Neural Voice** – Endpoint hosting | **$4.04 per hour per voice model** for hosting the custom voice endpoint [16] . This is charged for each hour the custom voice endpoint is running (billed per second). If the endpoint is stopped (suspended) when not in use, hosting charges pause [19] . One custom voice counts as one model/endpoint. |
| **Custom Neural Voice** – Long Audio | **$100 per 1 million characters** (same rate as other voices) when using long-form batch synthesis with a custom voice [15] . |
| **Other services** – Speech Translation, Speaker ID, etc. | *Speech translation:* ~$2.50 per audio hour [20] . <br>*Speaker Recognition:* ~$5 per 1,000 identification/verification transactions [20] . <br>*Voice Profile Storage:* ~$0.20 per 1,000 profiles per day (with a free allowance of 10,000 profiles per month) [20] . |

*Pricing notes:* Azure's **Free Tier** offers limited free usage each month: e.g. 5 audio hours of standard STT, **0.5 million characters of neural TTS free per month**, etc. [21] . This allows developers to experiment at no cost up to those limits. Beyond that, the above pay-as-you-go rates apply. All billing is typically prorated – for instance, STT is billed in one-second increments [22] and custom endpoint hosting is billed per second of runtime [23] [19] . Text input for TTS is counted in characters (with 1 million characters roughly equivalent to about 700–800 pages of text). As shown, **custom voices carry additional costs** (training and higher usage fees) compared to standard voices, reflecting the specialized infrastructure and processing involved. Azure also requires an application and approval to use custom voices (next section), so there may be an initial overhead in obtaining access.

**Table 1** above focuses on core pricing. For up-to-date and detailed pricing (which may change over time), Microsoft provides an official [Azure AI Speech pricing page [16] [24] ] and a pricing calculator tool. Prices here are in USD and standard rates; enterprise customers or high-volume users might get discounts via Azure's enterprise agreements. It's also important to manage usage (e.g. shutting down custom endpoints when not needed, to avoid continuous hosting charges [19] ).

## Prerequisites and Restrictions for Using Custom Neural Voice

Because of the potential for misuse of voice cloning technology, Microsoft has put strict **prerequisites and safeguards** in place for Custom Neural Voice. **CNV is a Limited Access feature**, meaning you **must apply and be approved by Microsoft** before you can create or deploy a custom voice model [25] . This approval process is part of Microsoft's Responsible AI commitment to ensure synthetic voices are used ethically and with proper consent.

**Application process:** To get access, customers must submit a **Limited Access registration form** detailing their intended use case for custom voice [26] . Only certain use cases are permitted – generally those that are non-malicious, such as assistive technologies, branded voice assistants, content creation with consent, etc. Impersonation of individuals without consent or deceptive uses are strictly prohibited. Microsoft reviews each application; only customers who are managed by a Microsoft account team or who clearly meet the criteria may be approved [27] . The use case(s) you propose in the application, if approved, define what you're allowed to use the custom voice for. Microsoft may require re-validation of the project periodically. In short, **not everyone can immediately use CNV** – you need to justify a responsible use case and agree to the terms.

**Consent and voice talent requirements:** If the custom voice is based on a voice actor or any person's recordings (including your own voice), you must obtain explicit written **consent from the voice talent** before creating the voice model [28] . In fact, as part of the process, Microsoft requires a *voice talent acknowledgement recording*: the voice actor must record a statement acknowledging the use of their voice in creating an AI model [29] . Azure uses speaker verification to compare this acknowledgement with the training audio to ensure the person consented is indeed the voice in the training data [29] . This safeguard helps prevent someone from covertly cloning a voice without the speaker's knowledge. Moreover, customers are contractually bound to only use the custom voice for the approved scenarios and to have proper rights to the voice recordings [28] .

**Responsible AI and usage guidelines:** All users of Custom Neural Voice must adhere to Microsoft's **Code of Conduct for Text-to-Speech** and the synthetic voice deployment guidelines [30] . Key points include: making sure people are **informed when they are listening to a synthetic voice** (transparency), not using the technology to deceive or impersonate in fraudulent ways, and including appropriate **disclosures** in applications. Microsoft provides guidance on how to disclose synthetic voice usage (e.g., a message like "This voice is AI-generated using a synthetic voice"). There are also recommendations for watermarking and detection – in fact, Azure's system automatically injects subtle digital **watermarks** into audio generated by custom voices and personal voice, which help in identifying that the audio is AI-synthesized [31] [32] . Microsoft continually improves these watermarks for robustness. The **limited access terms** further emphasize that users should not attempt to create voices of public figures or anyone without consent, and that any voice talent used has signed the required forms [28] .

In summary, to use Custom Neural Voice you need to: (1) **Apply and be approved** for access, (2) **Obtain consents** from any voice talent and follow the provided scripts for recordings, (3) **Abide by usage policies** (only use the voice for the approved purpose, ensure transparency to end-users, no harmful use), and (4) **Protect the voice data** (adhere to privacy and security requirements). These prerequisites protect individuals' rights and help prevent the creation of deceptive "deepfake" audio [25] . They also mean that launching a custom voice project may involve lead time for approval and preparation, so plan accordingly. For those who just want to experiment, Azure's **CNV Lite** or **Personal Voice demo** in Speech Studio can be tried (in preview) to create a rudimentary custom voice with limited data [5] – but even that requires registration to access the demo in most cases [33] .

## Real-World Use Cases and Case Studies

Azure's custom voice capabilities have been used in a variety of innovative applications, ranging from branded virtual assistants to accessibility solutions. Below we highlight several real-world examples and success stories that demonstrate how Custom Neural Voice and related speech services are applied:

*Progressive's "Flo" virtual assistant gained a friendly, recognizable voice thanks to a Custom Neural Voice model cloned from the actress who portrays Flo [34] . This allowed Progressive Insurance to use the same persona in voice interactions as in their ads, creating a consistent branded experience on smart speakers and chatbots [35] .*

- **Virtual Assistants and Brand Voices:** One of the early adopters of Custom Neural Voice was **Progressive Insurance**, which created a custom voice for their virtual agent "Flo." Progressive had a popular chatbot of Flo on text channels and wanted to extend her persona to voice platforms. Using CNV, they generated a synthetic voice that sounds like Flo (voiced by actress Stephanie Courtney) to power the Flo chatbot on smart speakers and IVR systems [34] [35] . This gave customers the same playful, no-nonsense voice they recognize from TV ads when interacting with Progressive's voice assistant. Another example is **AT&T's Bugs Bunny experience** – to showcase 5G, AT&T built an augmented reality app where Bugs Bunny talks to customers in real time. They worked with Warner Bros. and Azure to create Bugs's iconic cartoon voice as a custom neural voice. Visitors could talk to an AR Bugs Bunny who greets them by name and responds with a perfect Bugs Bunny voice, enabled by the Azure custom voice model [36] [37] . These cases illustrate how companies use custom voices to **personify characters and brands** in a unique way.

- **Personal Voice Assistants:** With the introduction of Personal Voice, individual users can bring their own voice into their devices. For instance, the caller ID app **Truecaller** integrated Azure's personal voice in their **Truecaller Assistant**, a call-screening digital assistant. Truecaller's Assistant answers incoming calls for the user, converses with the caller to filter spam, and summarizes the call – and now it can do so *using the user's own voice*. When a Truecaller user enables this feature, they provide a short recording and the app creates a custom voice so that the assistant sounds like the user, adding familiarity and comfort for callers [38] [39] . According to Truecaller, this personalization makes the call-screening experience more engaging and transparent to the caller (it feels like the person is speaking to them, even though it's an AI) [40] [41] . This is a novel use of custom voice in the realm of personal productivity and phone applications.

- **Accessibility and Voice Preservation:** One of the most impactful use cases for custom voice is giving people their voice back. Microsoft MVP **Charles Elwood** created an app that uses Custom Neural Voice to help individuals who lost the ability to speak (due to conditions like throat cancer, ALS, or other diseases) by building a synthetic voice bank [42]. In one case, a former radio DJ named Chris lost his voice after larynx removal surgery. Because recordings of his voice existed from his radio days, Charles used those clips to train a custom neural voice model of Chris's voice [43] [44]. Now, Chris can type into an app and have it speak in *his own voice* – allowing him to communicate with family and even resume activities like announcing sports, which he loved [45]. His wife remarked that the technology "brought him home" because it restored a core part of his identity [45]. This highlights how CNV can be an **assistive tool** for voice banking: people at risk of losing their voice can record samples ahead of time and later use a synthetic voice to speak. It offers a level of personal connection that generic TTS voices cannot match. There are ongoing pilots (e.g. with the **Team Gleason** foundation for ALS) exploring this technology for patients, showing the promise of Azure's custom voice in accessibility.

- **Multilingual Translation and Dubbing:** Azure's speech services also enable cross-language scenarios – effectively "dubbing" a voice into other languages. A prime example is **Skype's TruVoice Translator** feature, which leverages Azure AI. When two people chat via Skype in different languages, Skype can live-translate the speech and play it back in the other person's language. Now, with Azure's Personal Voice models, the translated speech is generated in the *original speaker's voice*, preserving the speaker's vocal identity even while speaking a different language [46]. For instance, if an English speaker talks to a French speaker, the French listener will hear a French translation spoken in the English user's own voice (synthetically reproduced). This dramatically improves the naturalness of translated calls compared to using a generic voice. It's essentially real-time AI dubbing of a conversation, and it uses custom neural voice technology under the hood [46]. Beyond live translation, the cross-language capability of CNV means content creators could dub videos or lectures into multiple languages using the same custom voice, maintaining a familiar voice across languages. This has potential in media localization – e.g., a narrator could have their voice "cloned" and used to narrate the same documentary in several languages, without needing the original person to speak those languages.

- **Education and Entertainment: Duolingo**, the popular language-learning app, leveraged Custom Neural Voice to create distinct voices for its cast of characters in the app [47] [48]. Duolingo introduced a lineup of quirky characters (Lily, Junior, and others) who speak sentences to the learner during lessons. Using CNV, Duolingo worked with voice actors to give each character a unique synthetic voice that can speak multiple languages. For example, they launched "Lily's" voice in both English and Spanish, and "Junior's" voice in English, with plans for voices in French, German, Japanese, etc., so that as you use Duolingo, the characters consistently maintain their personas across languages [48]. This personalization makes the learning experience more engaging and exposes learners to a variety of voice styles and accents [49]. On the entertainment side, beyond the Bugs Bunny AR example, we see custom voices being used in games and interactive content to voice fictional characters. Because CNV can create *composite voices* (combining qualities of different voices to make a new fictional voice) [50], game studios or animation companies can craft voices for characters that don't correspond to a single human – expanding creative possibilities while still sounding natural.

- **Content Creation (Audiobooks, Videos, Media):** Custom voices are beginning to be used in content production workflows. A notable example is **Wondershare**, which makes the popular video editing software Filmora. Wondershare is integrating Azure's personal/custom voice capability into their tools so that creators can generate voice-overs in their own voice for videos [51] [52] . Imagine a YouTuber typing out a script and the software producing the narration in the YouTuber's voice without them having to record it manually – this can save time and allow easy editing of the narration by just changing the script. It also helps maintain consistency if the creator is unavailable to record or needs to produce content in multiple languages. We're likely to see more of these content-creation applications, although Microsoft currently restricts use of custom voices for synthetic media to ensure it's done responsibly (they require those use cases to meet the limited access criteria and sometimes have additional oversight [53] ).

These case studies underscore the versatility of Azure's custom voice technology. From enhancing brand identity (e.g., Flo and Bugs Bunny) to **empowering individuals** (restoring someone's voice, personal call assistants) and enabling new **cross-language experiences**, custom neural voices are opening new horizons. Microsoft's approach, however, balances innovation with responsibility – all these deployments were done with consent and oversight. For instance, Progressive's Flo voice was created with the actress's involvement and used in a context where users know they're talking to a bot; Duolingo's character voices were created to augment learning, not to deceive; and the accessibility cases are done with deep respect for the individuals involved.

## Conclusion

Azure's speech AI offerings provide a comprehensive toolkit for voice interaction – from high-quality text-to-speech and speech recognition to the crown jewel: Custom Neural Voice for building bespoke voices. Using custom voices, organizations and developers can create more engaging and personalized experiences, whether it's a **virtual assistant with a distinctive brand voice**, an **accessible communication device that speaks in a patient's own voice**, or a **multilingual app that retains character voices across languages**. Microsoft has invested heavily in the realism of these voices, leveraging neural networks to capture the nuances of human speech [8] . At the same time, Azure imposes strict **responsible AI practices** – ensuring consent, security, and transparency – to mitigate risks associated with voice cloning [25] [29] .

For those interested in using a custom voice: start by exploring the standard Speech service capabilities (TTS and STT) to see if they meet your needs. Azure's prebuilt neural voices are continually improving and might already provide the style you want. If a truly unique voice is needed, you can pilot with *Custom Neural Voice Lite* (record ~20-50 sentences to create a test voice) [5] or the Personal Voice demo, to evaluate the technology. Then, prepare a solid use case and apply for access to Custom Neural Voice. Ensure you line up all the requirements – voice talent agreements, scripts, and a plan for responsible use. Once approved, Azure provides the tools (Speech Studio, SDKs, and documentation) to guide you through dataset preparation, model training, and deployment.

Azure's custom voice technology has already seen success in industry, education, accessibility, and entertainment. As speech synthesis continues to advance, we can expect even more natural and versatile voices. With Azure's services, developers can be at the forefront of this innovation – giving applications a unique voice while upholding the trust and ethics that users expect.

**Sources:** Microsoft Azure Cognitive Services documentation and pricing pages; Microsoft Tech Community and Azure Blog posts; Microsoft case studies and news releases; and other reputable tech sources as cited throughout. The information is current as of 2024-2025, but readers should refer to official Azure docs for the latest details.

---

[1] [2] [3] [4] [9] [30] Custom neural voice overview - Speech service - Azure AI services | Microsoft Learn
https://learn.microsoft.com/en-us/azure/ai-services/speech-service/custom-neural-voice

[5] [25] [26] [27] [28] [33] Limited Access - Azure AI services | Microsoft Learn
https://learn.microsoft.com/en-us/legal/cognitive-services/speech-service/text-to-speech/limited-access?context=%2Fazure%2Fai-services%2Fspeech-service%2Fcontext%2Fcontext

[6] [7] [31] [32] [38] [39] [40] [41] [46] [51] [52] [53] Create personalized voices with Azure AI Speech
https://techcommunity.microsoft.com/blog/azure-ai-services-blog/create-personalized-voices-with-azure-ai-speech/4147073

[8] [29] [34] [35] [36] [37] [47] [48] [49] [50] Are you talking to me? Azure AI brings iconic characters to life with Custom Neural Voice - Source
https://news.microsoft.com/source/features/ai/custom-neural-voice-ga/

[10] [11] [12] [13] [22] [23] Azure AI Speech Pricing | Microsoft Azure
https://azure.microsoft.com/en-us/pricing/details/cognitive-services/speech-services/

[14] [15] [20] [21] [24] Microsoft Azure Terms of Service | Speechify
https://speechify.com/blog/microsoft-azure-terms-service/?srsltid=AfmBOop2oGxOD9q4MrgKRDN-dxXyLWg6fwsedMdhD9VlzBUN4iIc9nuV

[16] [17] [18] [19] Custom neural voice pricing - Microsoft Q&A
https://learn.microsoft.com/en-us/answers/questions/1346760/custom-neural-voice-pricing

[42] [43] [44] [45] Tech for Good: Giving a Voice to the Voiceless | Microsoft Community Hub
https://techcommunity.microsoft.com/blog/mvp-blog/tech-for-good-giving-a-voice-to-the-voiceless/4217303