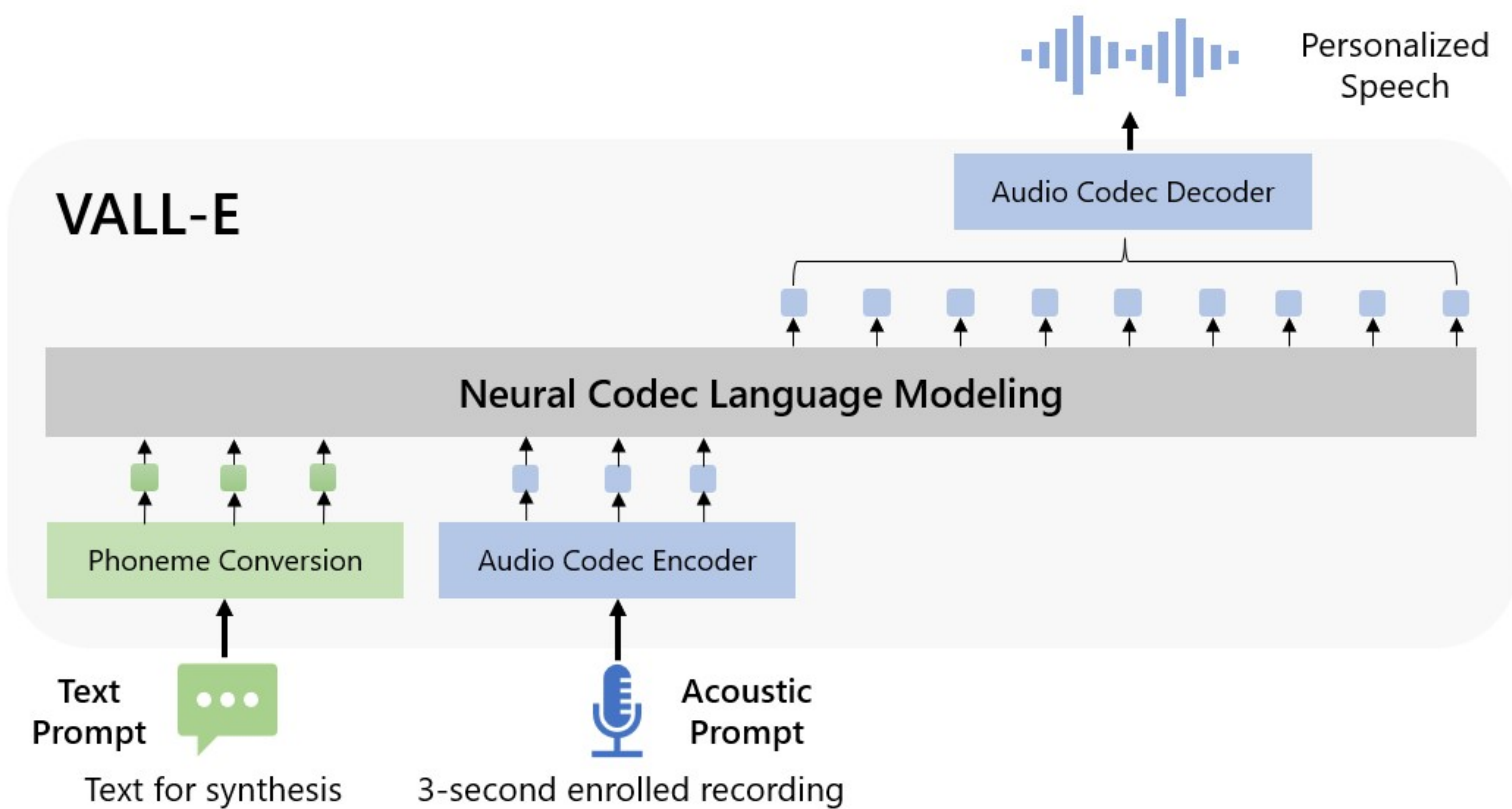# The state of the art (2 of 2)

- Class slides

# Orientation

- <u>Large speech language models</u>
  - VALL-E

- Tasks beyond Text-To-Speech

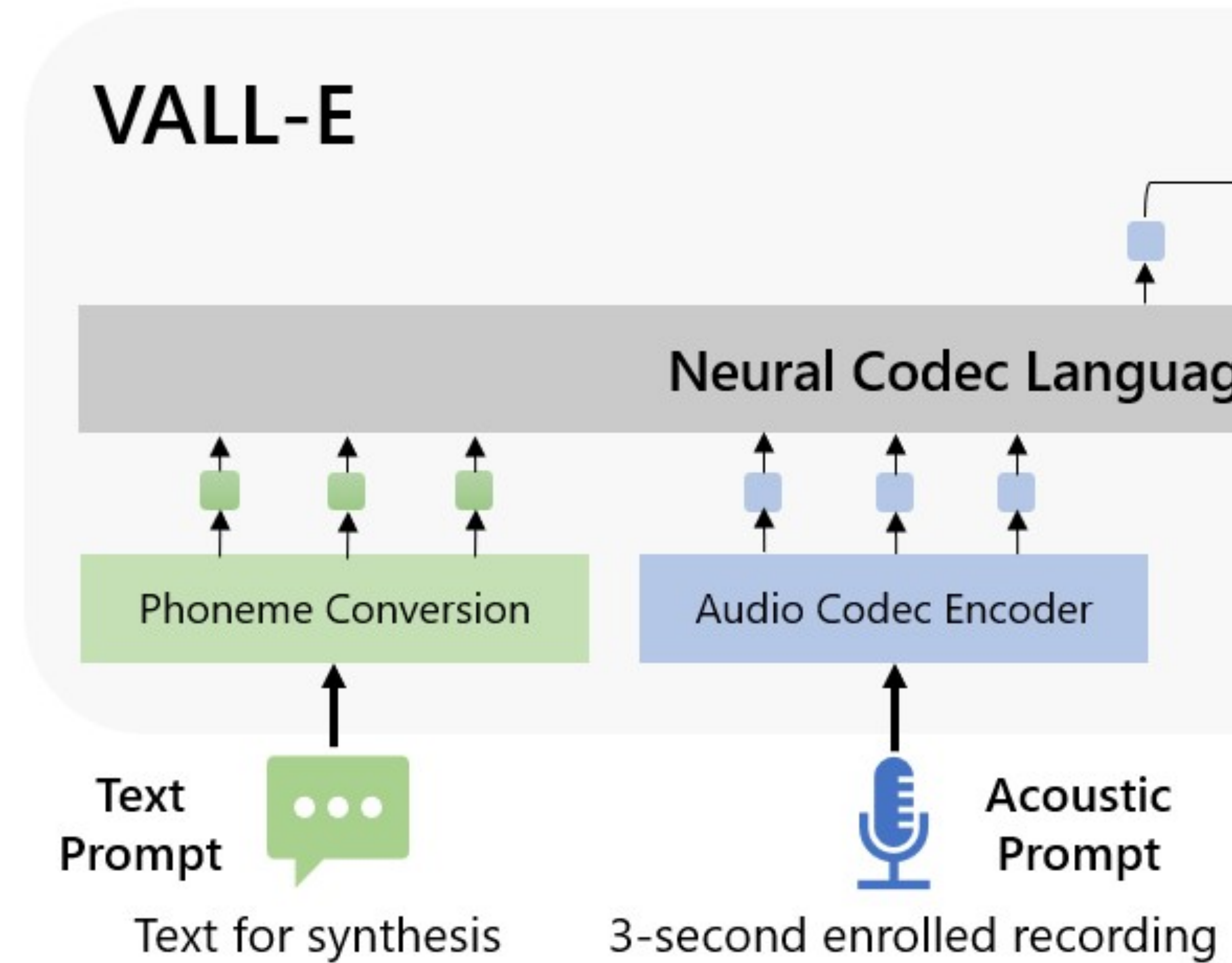- Current & future trends

# VALL-E



Personalized Speech

Audio Codec Decoder

Neural Codec Language Modeling

Phoneme Conversion

Audio Codec Encoder

**Text Prompt**

Text for synthesis

**Acoustic Prompt**

3-second enrolled recording

# VALL-E

Table 1: A comparison between VALL-E and current cascaded TTS systems.

|  | **Current Systems** | **VALL-E** |
| --- | --- | --- |
| Intermediate representation | mel spectrogram | audio codec code |
| Objective function | continuous signal regression | language model |
| Training data | $\leq$ 600 hours | 60K hours |
| In-context learning | ✗ | ✓ |

# How can we combine text and speech into a single sequence ?



VALL-E

Neural Codec Languag[e]

Phoneme Conversion    Audio Codec Encoder

**Text Prompt** — Text for synthesis

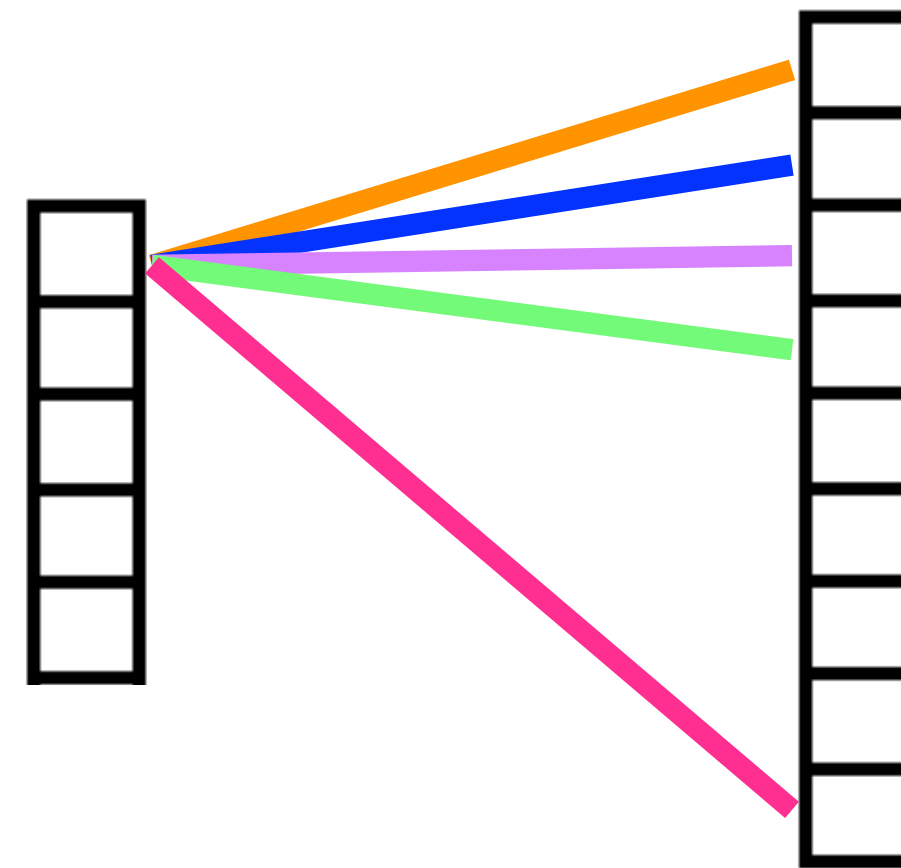**Acoustic Prompt** — 3-second enrolled recording

# Inputting a one-hot vector into the model: **embedding**
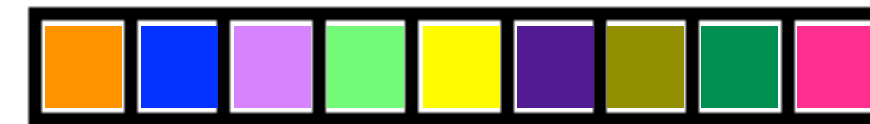
# Inputting a one-hot vector into the model: **embedding**

# Inputting a one-hot vector into the model: **embedding**

# Inputting a ~~one hot vector~~ symbol into the model: **embedding table**
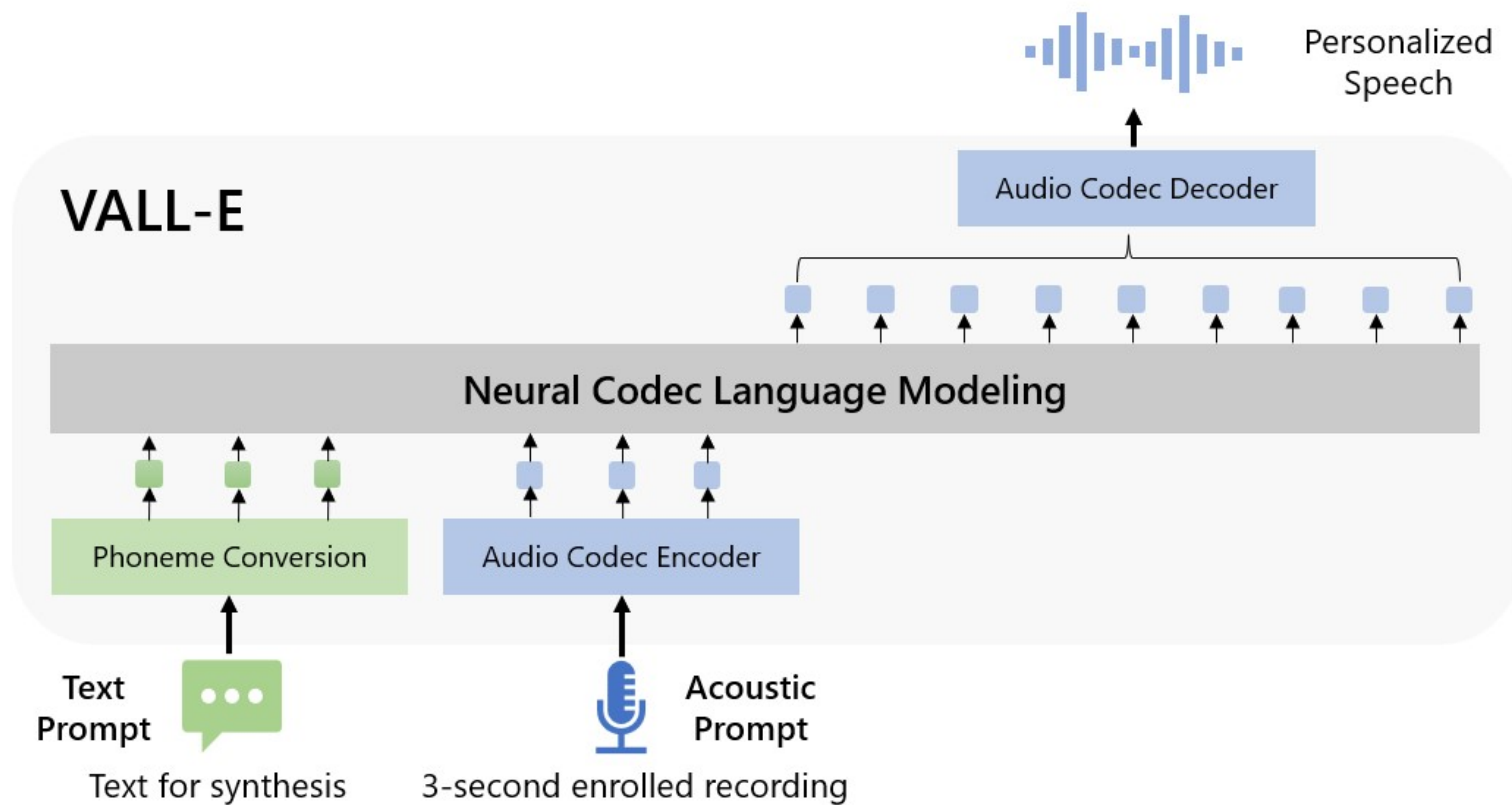
# How can we combine two different types of symbol into a single sequence ?
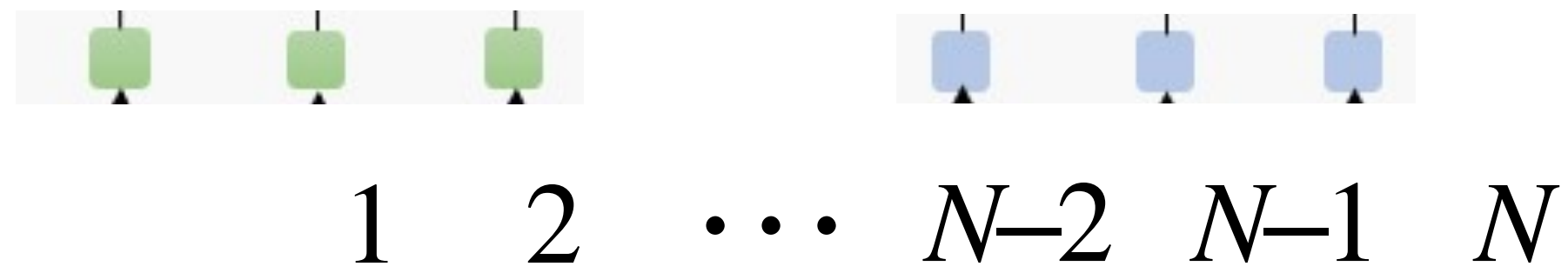
Option 1: a single embedding table

Option 2: separate embedding tables
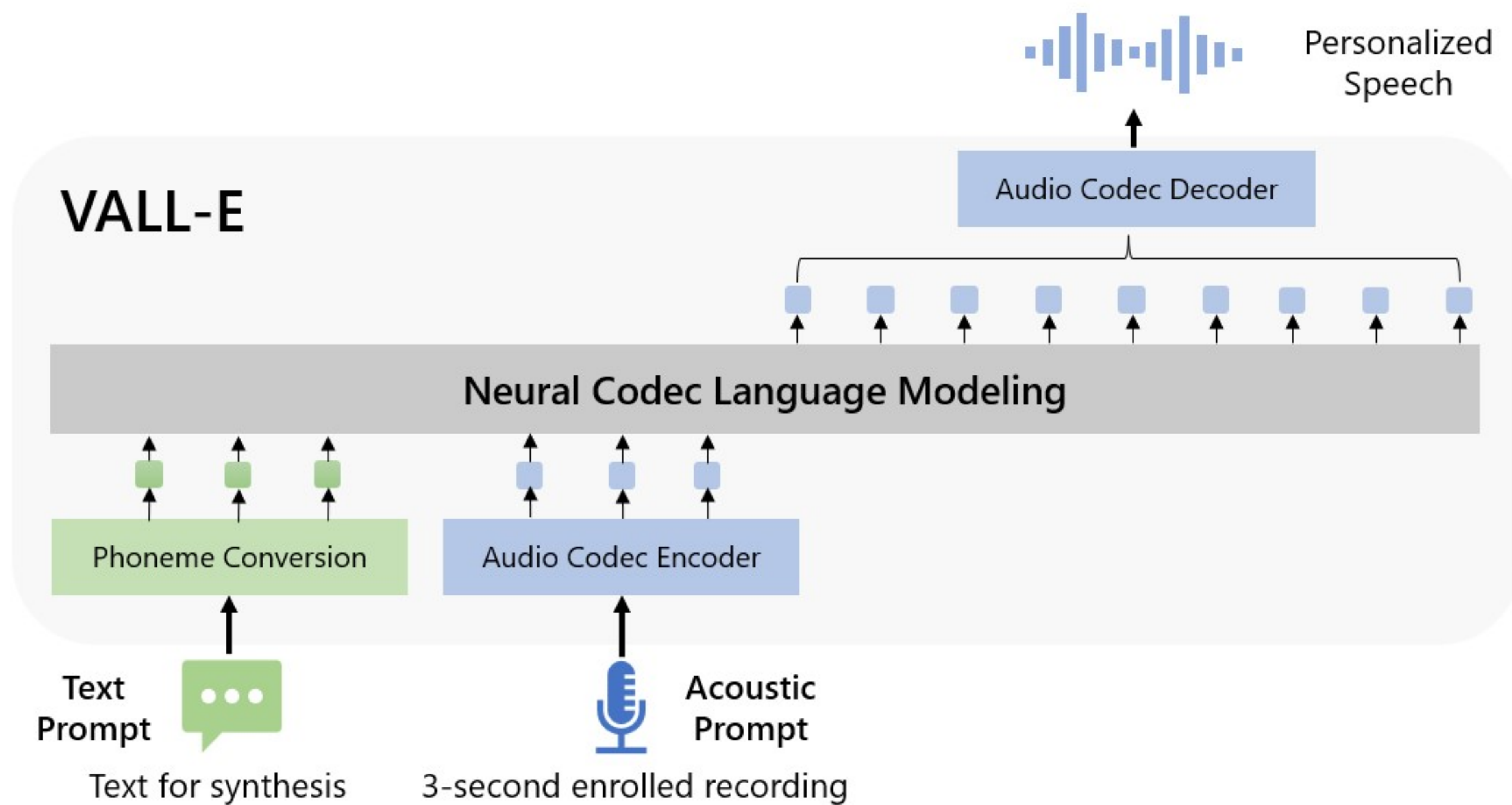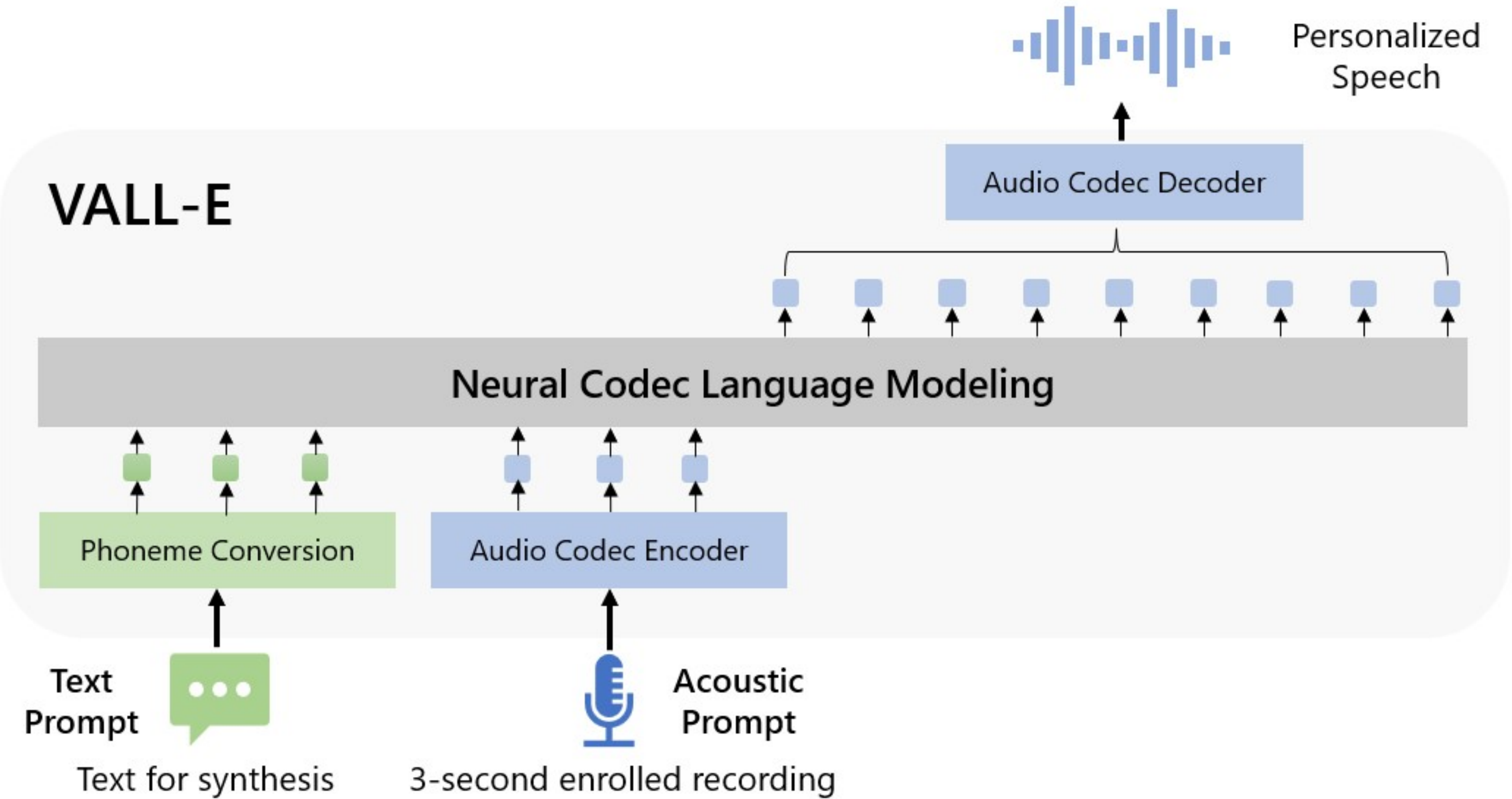
# Language modelling

# Language modelling

$$P(w_N \mid w_1, w_2, \ldots w_{N-1})$$

1    2    $\cdots$    N–2    N–1    N

# In-context learning (via prompting)

# Zero shot

# Orientation

- Large speech language models
  - VALL-E

- <u>Tasks beyond Text-To-Speech</u>

- Current & future trends

- Controllable TTS

- Voice conversion

- Prosody transfer

- Speech editing

- Speech translation

- ...etc

# Orientation

- Large speech language models
  - VALL-E

- Tasks beyond Text-To-Speech

- <u>Current & future trends</u>

- Larger models, larger data
- Pre-training
  - open models used as starting point by other researchers
  - fine-tuning and/or prompting
- Multi-task models
  - speech
  - music
  - "general audio"

# What next?

- Today's ''state-of-the-art'' will not last

- But understanding the history of TTS will help us understand what comes next

- Read the literature