

The state of the art (2 of 2)

- Class slides

Orientation

- Large speech language models
 - VALL-E slides from last week (recap)
- Tasks beyond Text-To-Speech
- Current & future trends



Orientation

- Large speech language models
- VALL-E
- Tasks beyond Text-To-Speech
- Current & future trends



Orientation

- Large speech language models
 - VALL-E
- Tasks beyond Text-To-Speech
- Current & future trends



Orientation

- Large speech language models
 - VALL-E
- Tasks beyond Text-To-Speech

- Current & future trends

Generation tasks

- Controllable TTS
 - “zero shot”
 - natural language description
- Speech editing
- Voice Conversion (VC)
- Voice privacy

Classification tasks

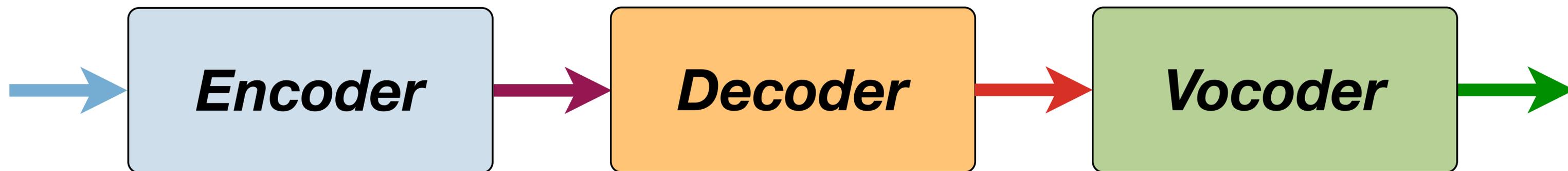
- Automatic Speaker Verification (ASV)
- “Anti-spoofing” / deepfake detection
- Source speaker tracing

Generation tasks

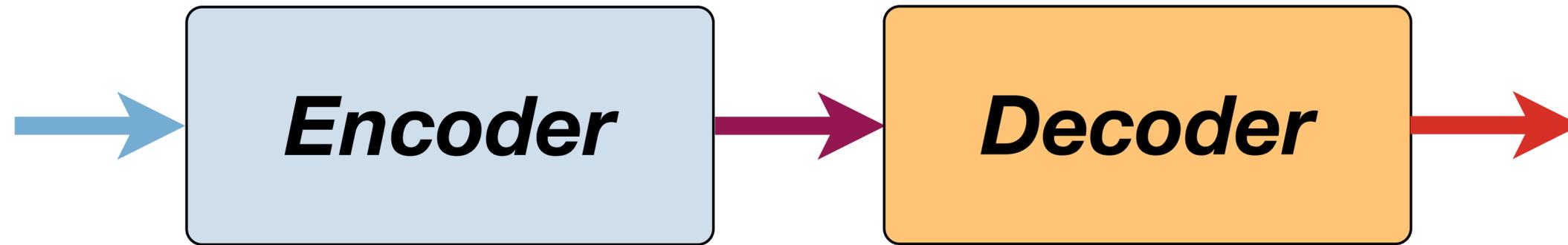
- Controllable TTS
 - “zero shot”
 - natural language description
- Speech editing

- Voice Conversion (VC)
- Voice privacy

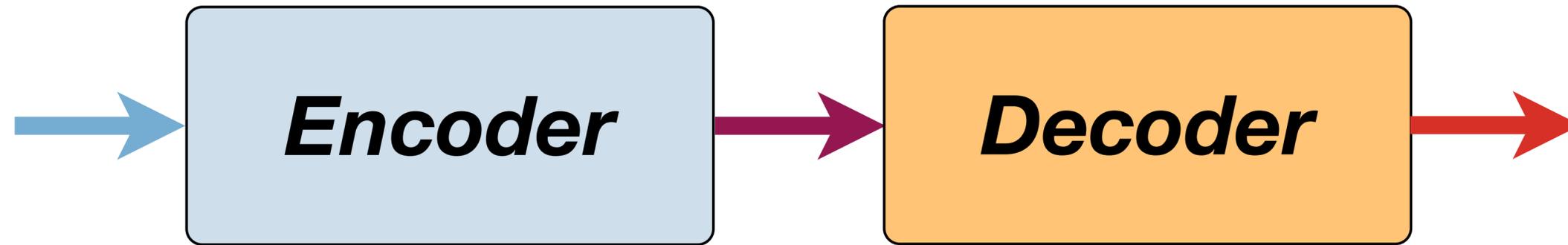
Voice Conversion



Voice Conversion - easy to train with parallel data, if you have some...



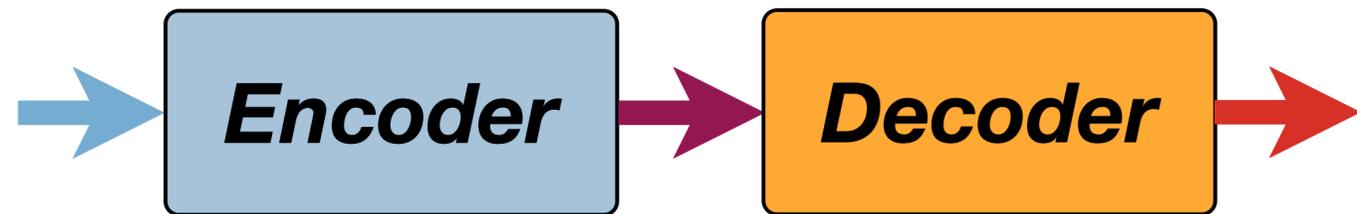
Voice Conversion without parallel data



Early systems:

1. create pseudo-parallel data,
2. train the system in the same way as for parallel data.

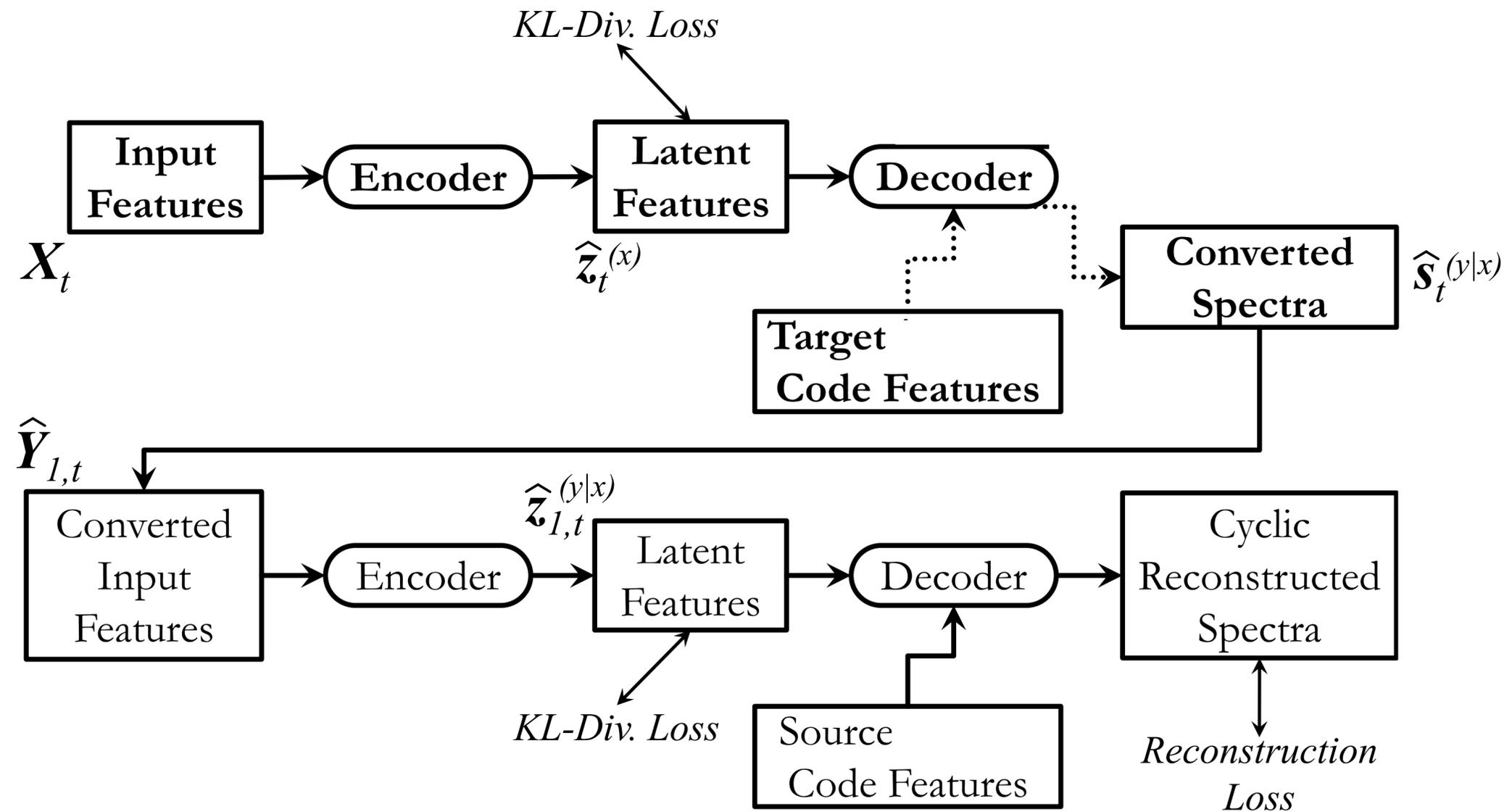
Voice Conversion without parallel data: **cycle-based approach**



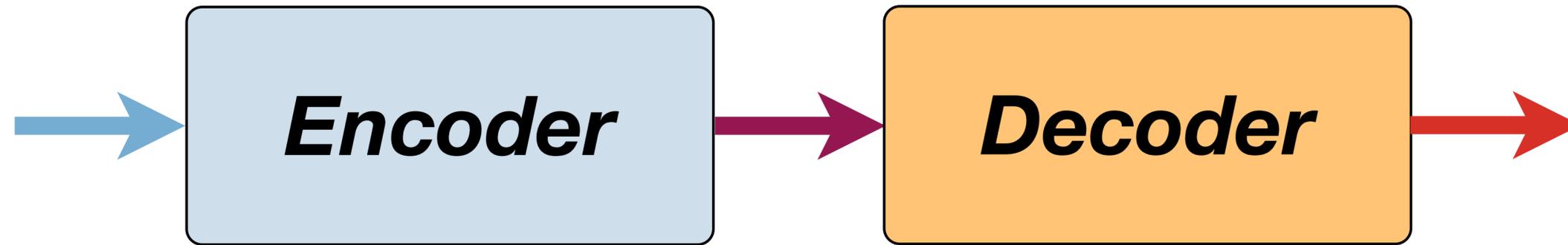
Cycle approach:

1. convert source-to-target; *cannot measure loss*
2. convert target back to source; now measure the loss.

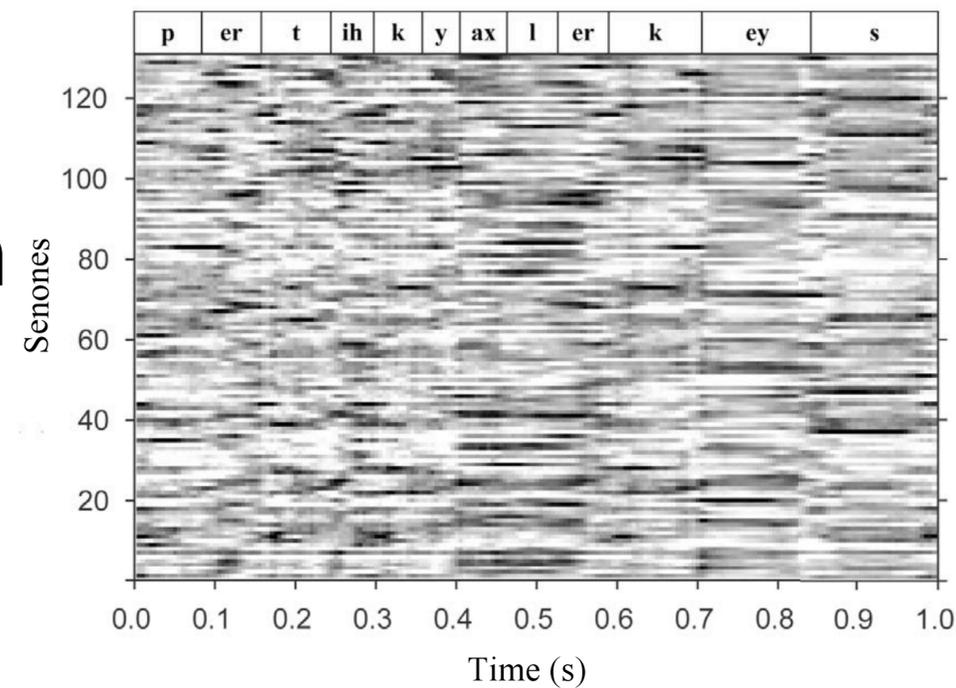
Cycle-based approach: source-target-source & source-source



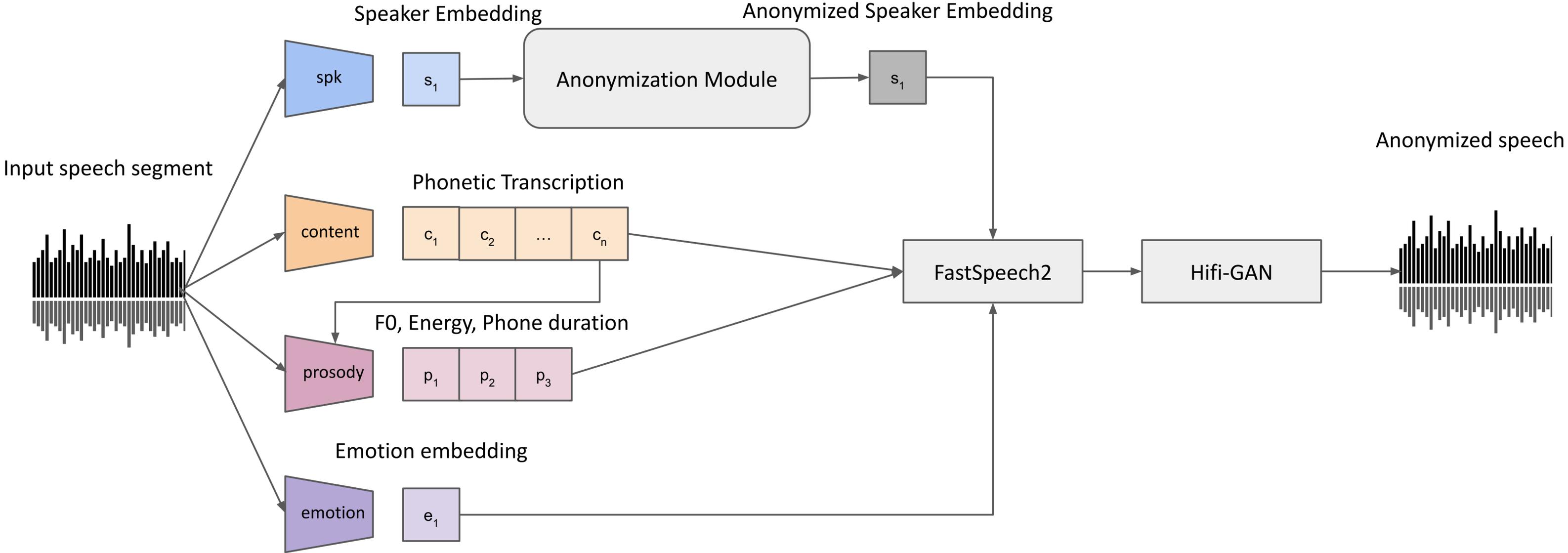
Voice Conversion without parallel data: **ASR+TTS** approach



Phonetic Posteriorgram
(PPG)



Voice privacy



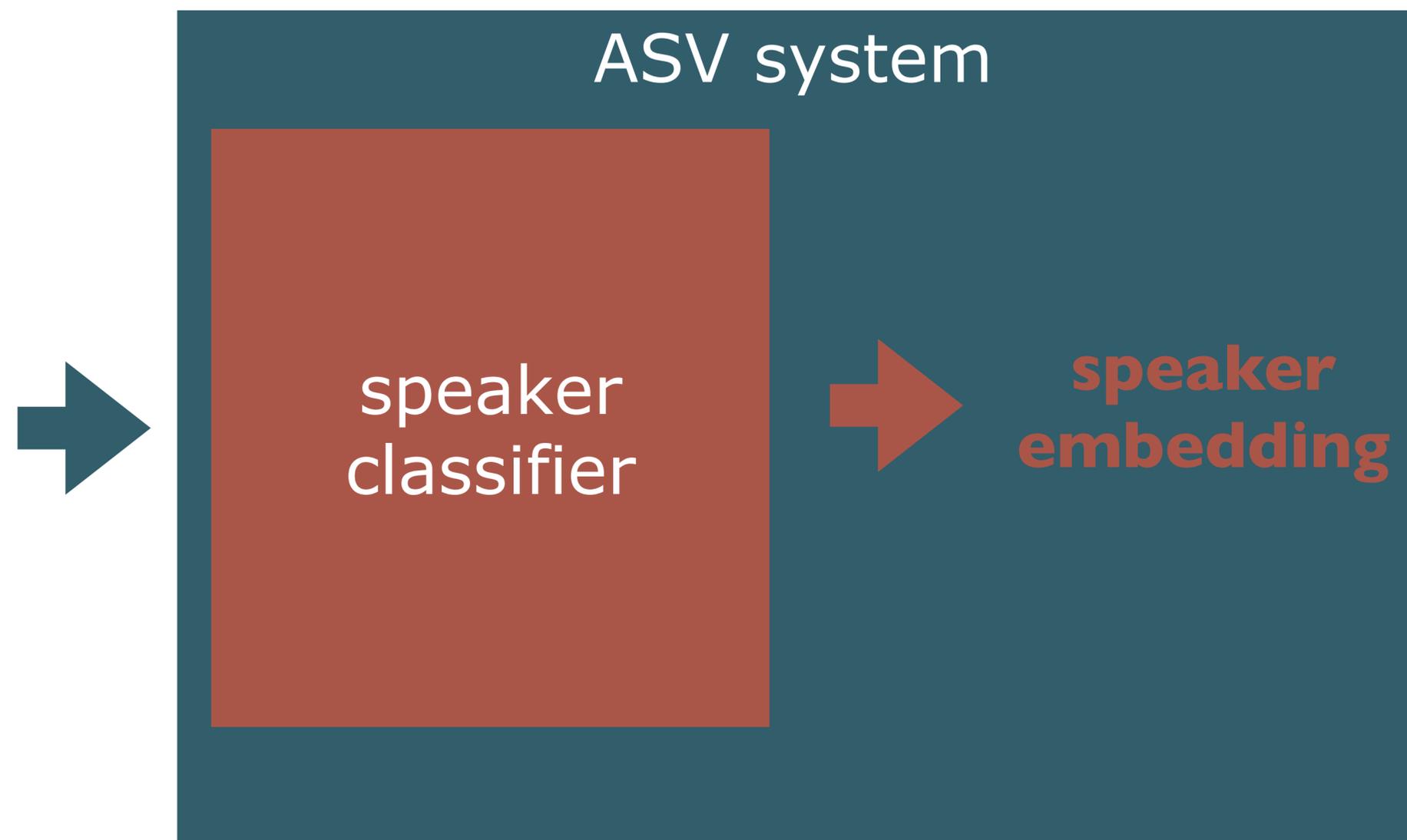
Classification tasks

- Automatic Speaker Verification (ASV)
- “anti-spoofing”
 - deepfake detection
- Source speaker tracing

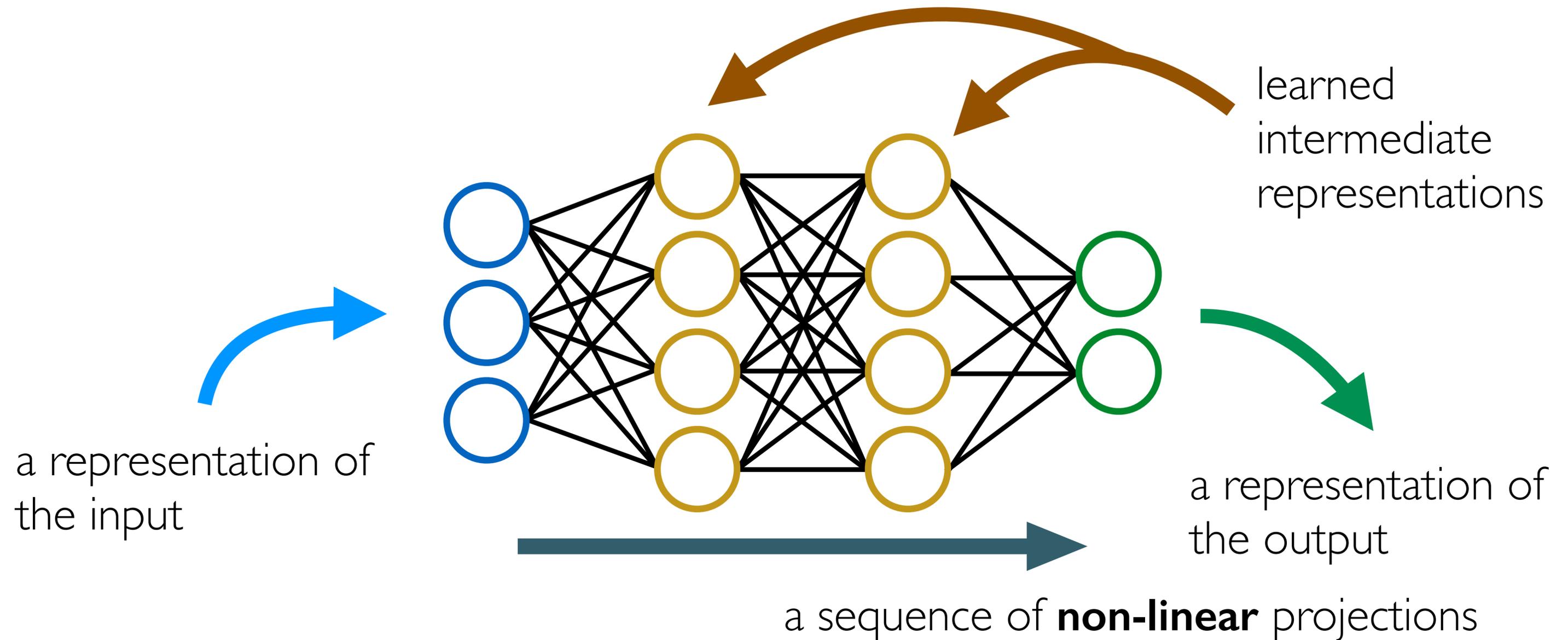
Automatic Speaker Verification (ASV)



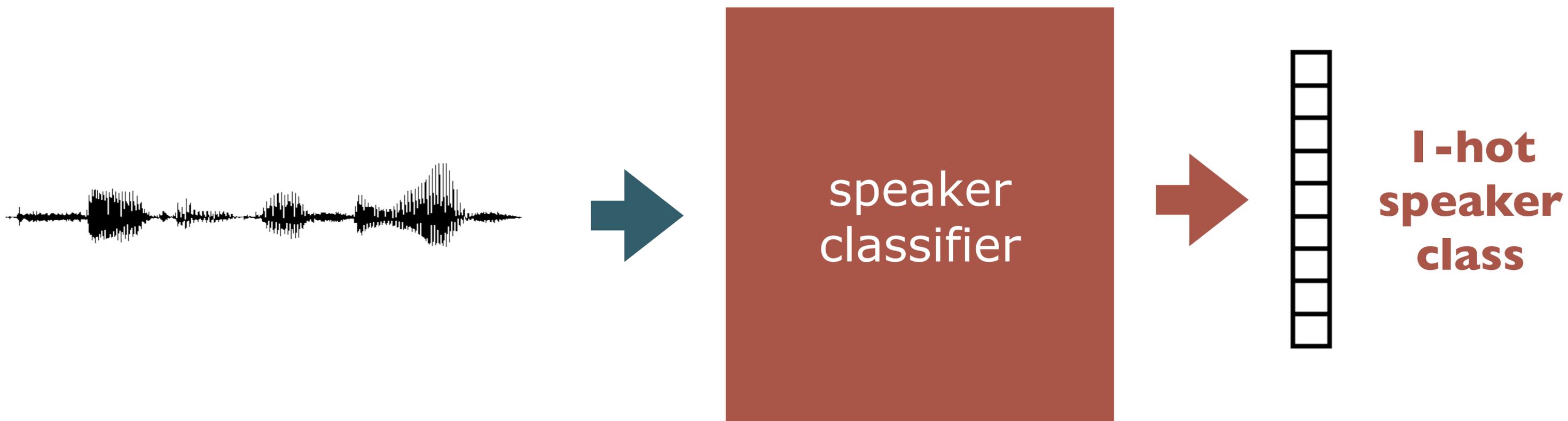
Automatic Speaker Verification (ASV)



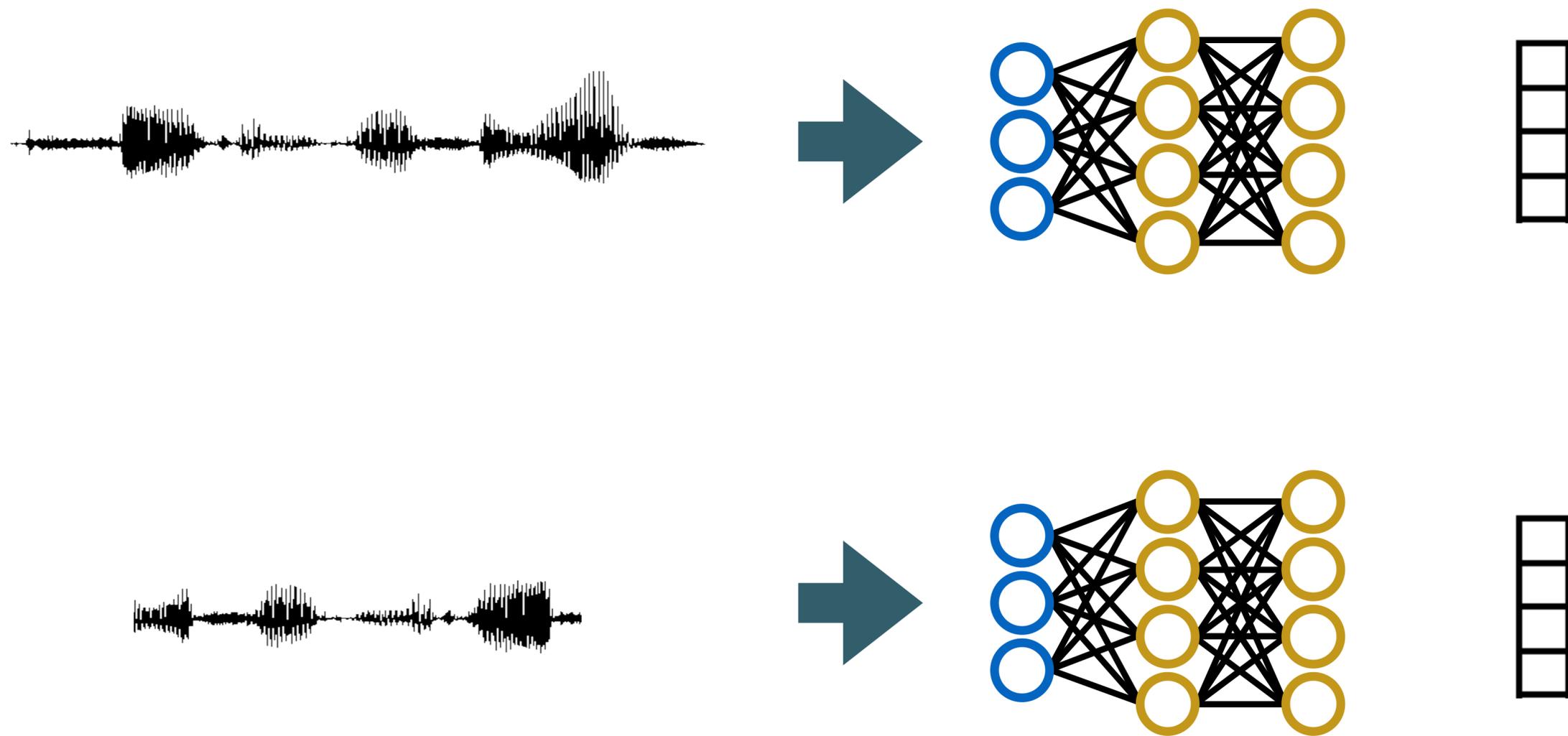
RECAP! What are all those layers for? Learning **representations**!



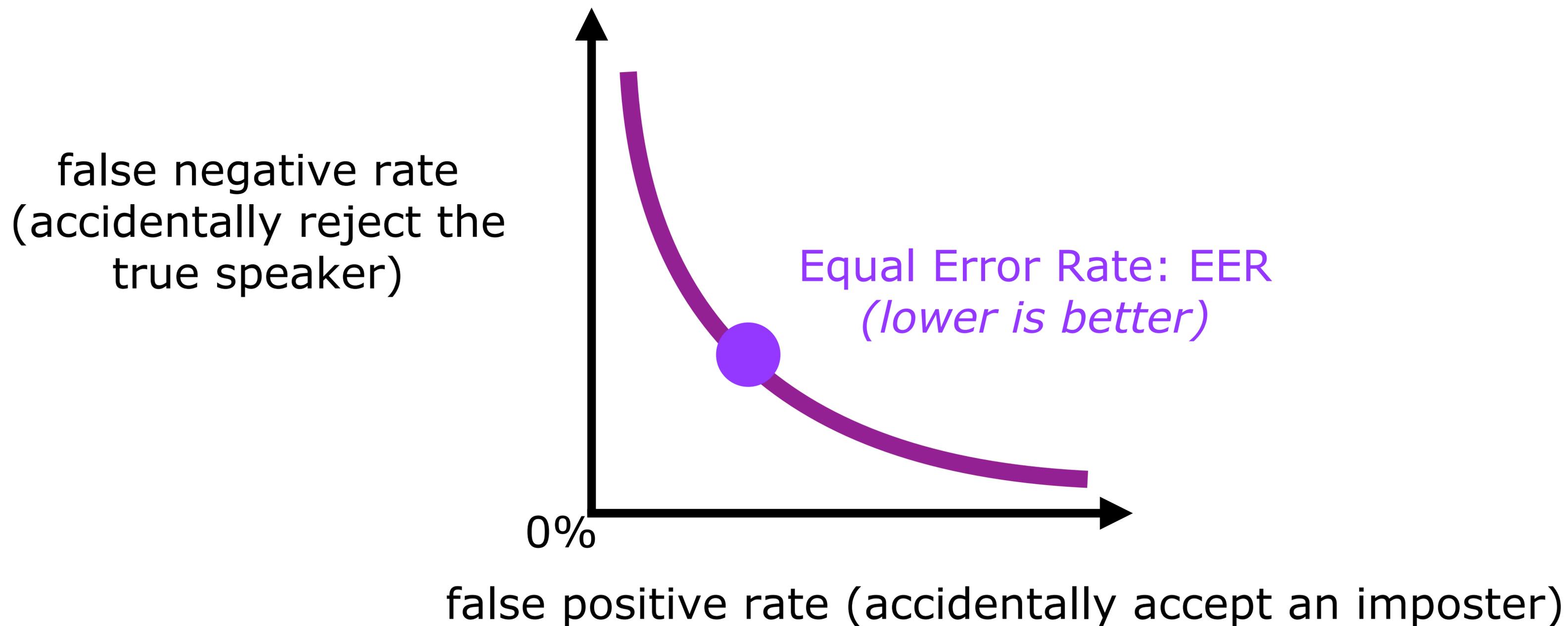
Speaker embedding



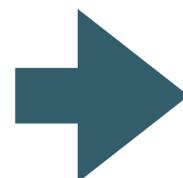
Speaker verification using speaker embeddings



Automatic Speaker Verification: error behaviour



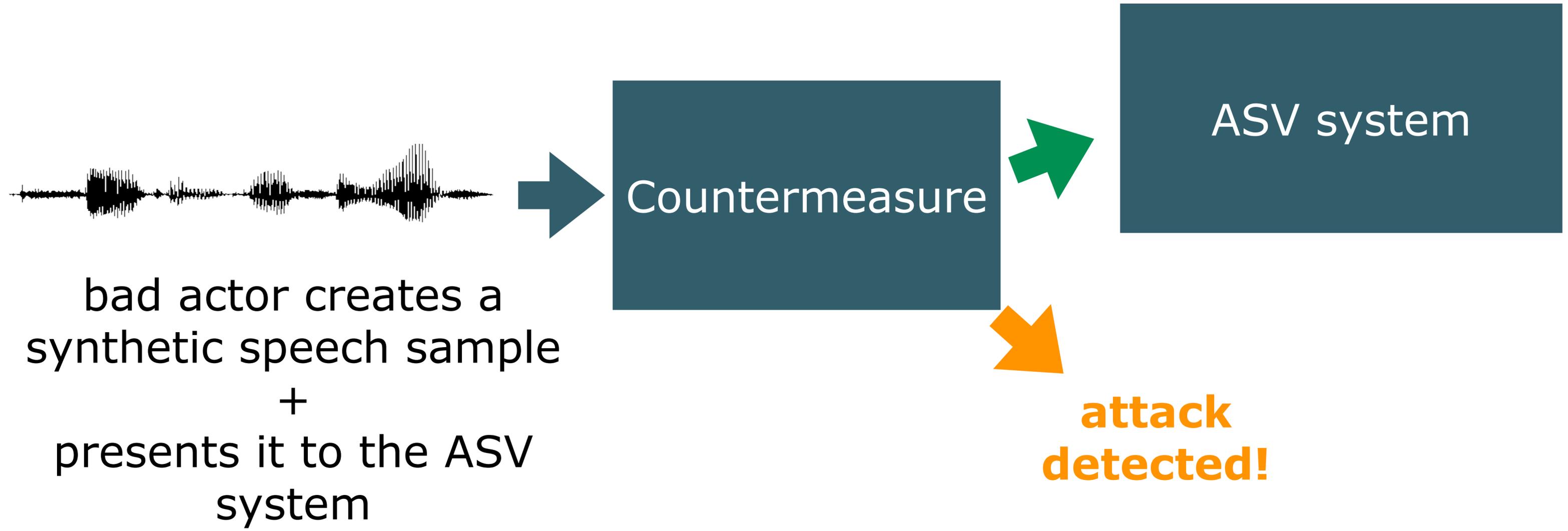
“spoofing” attacks on Automatic Speaker Verification



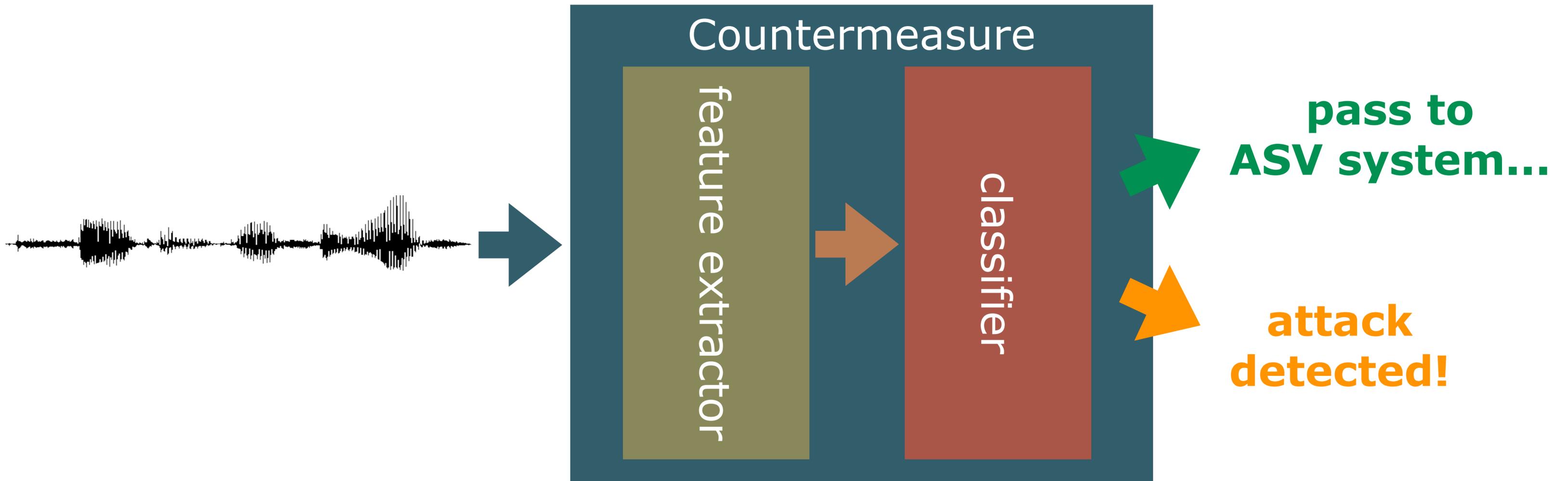
ASV system

bad actor creates a
synthetic speech sample
+
presents it to the ASV
system

Countermeasures against “spoofing” attacks on ASV

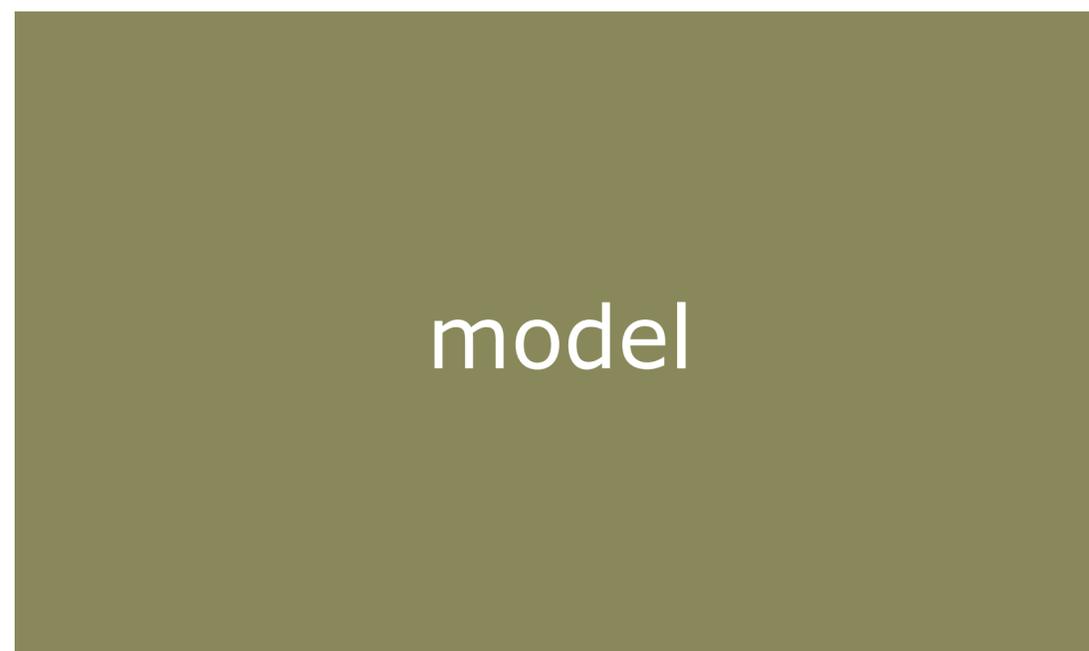


State-of-the-Art countermeasures: trained from data



Self-Supervised Learning (SSL)

“the cat sat on the mat”

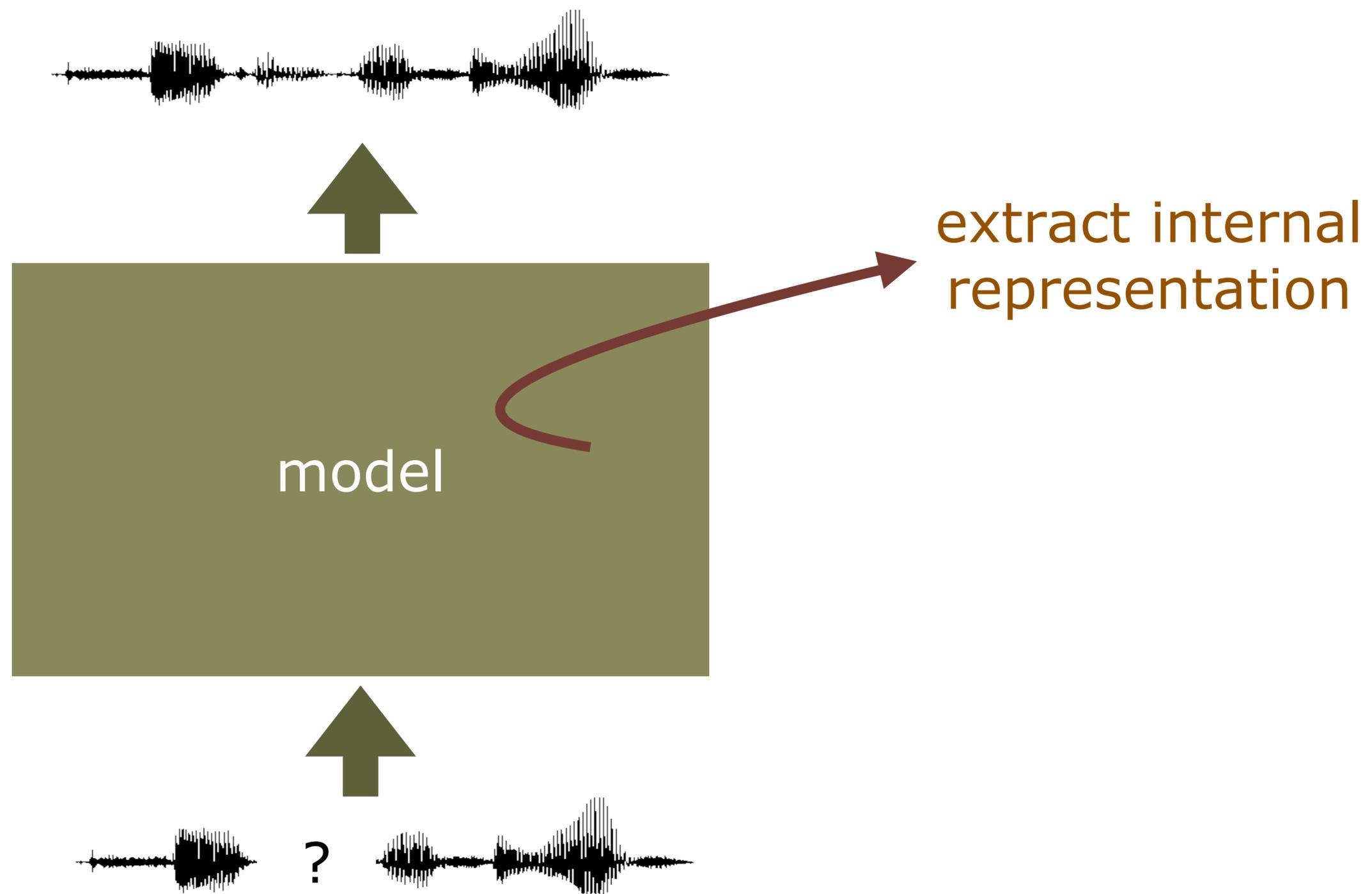


model

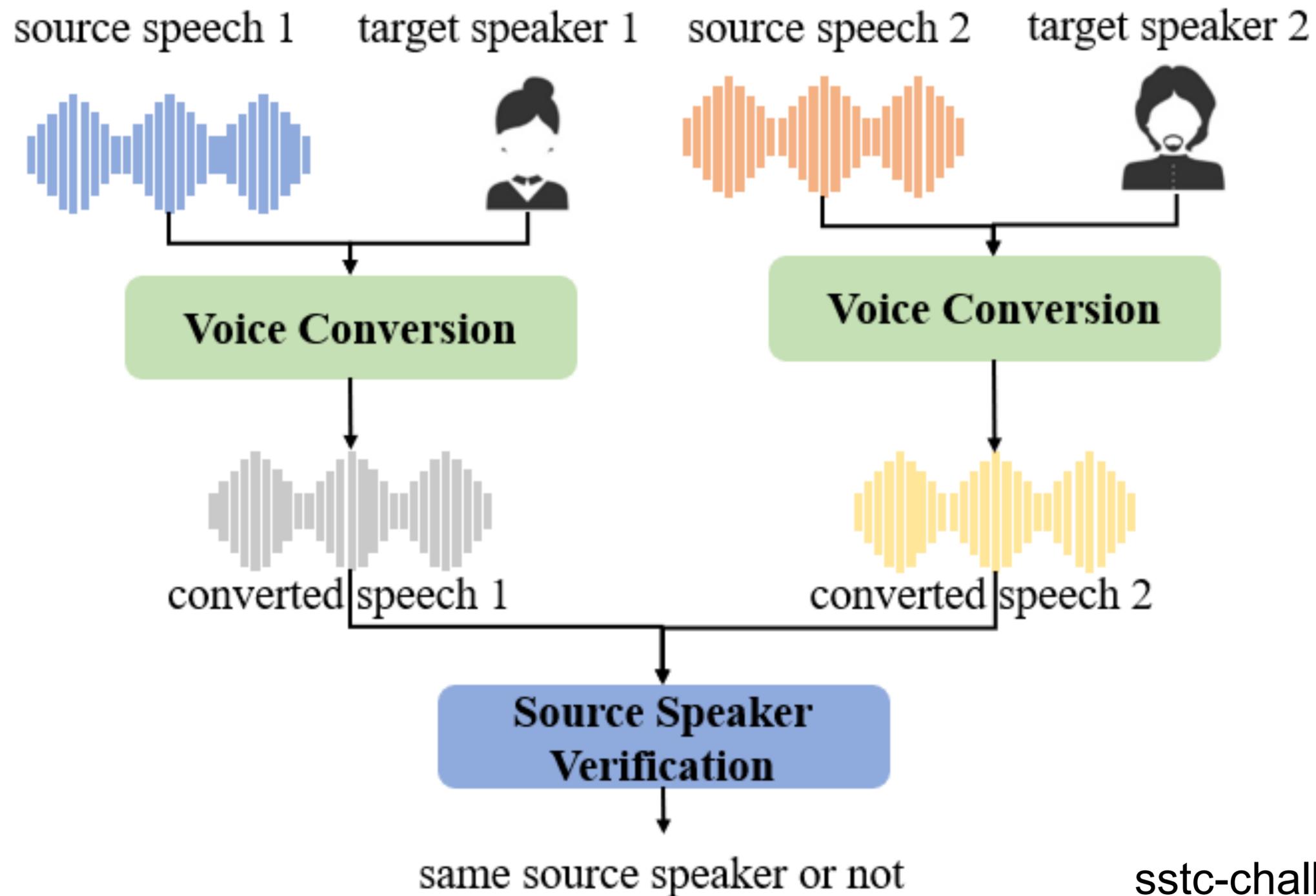


“the ? sat on the mat”

Self-Supervised Learning (SSL)



Source speaker tracing (here, just *verification*)



Orientation

- Large speech language models
 - VALL-E
- Tasks beyond Text-To-Speech
- Current & future trends



Orientation

- Large speech language models
 - VALL-E
- Tasks beyond Text-To-Speech
- Current & future trends
- 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

