

Speaking and Singing Voice Alignment with Deep Neural Networks



***Yann
Teytaut***



**Dr. Guillaume
Doras**



**Dr. Axel
Roebel**

Analysis/Synthesis – IRCAM (STMS UMR 9912) – Paris, France



**SORBONNE
UNIVERSITÉ**



**MINISTÈRE
DE LA CULTURE**

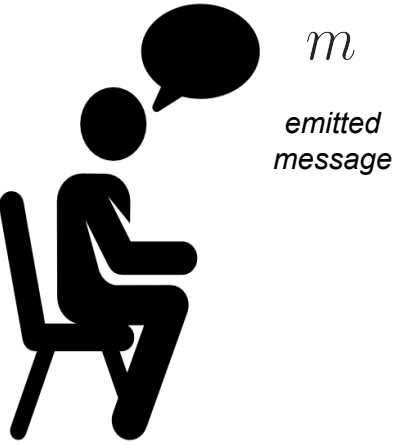
*Liberté
Égalité
Fraternité*



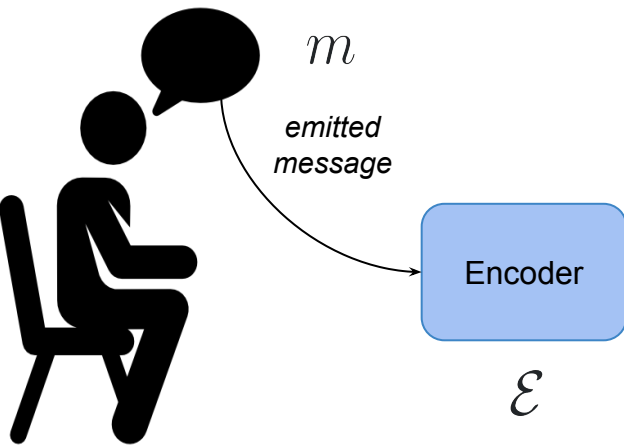
. Human communication



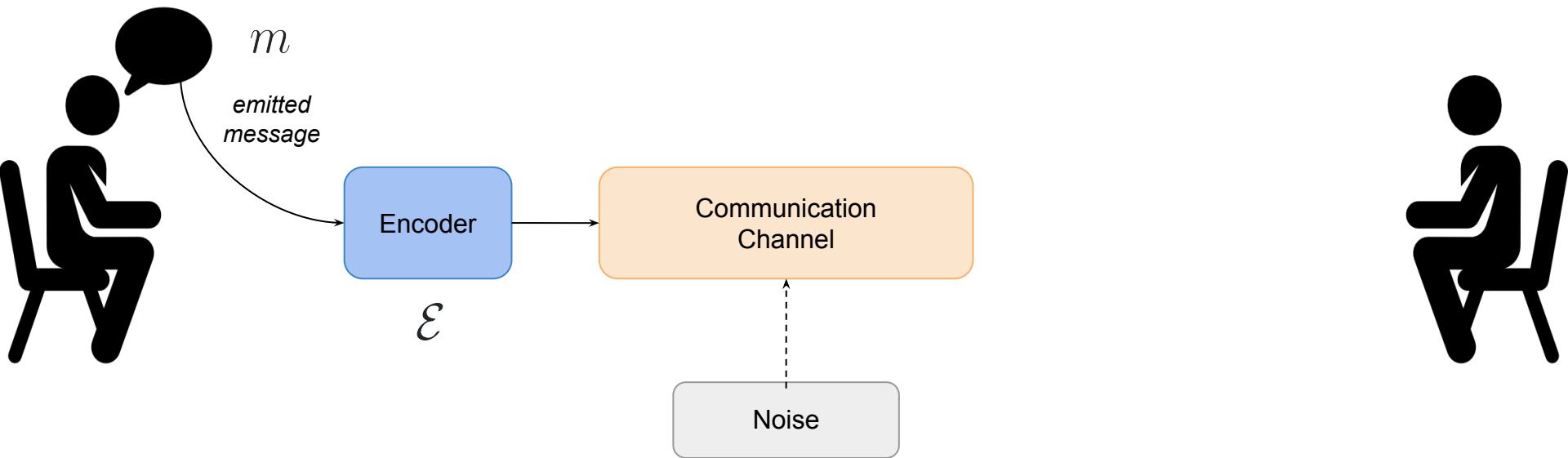
. Human communication



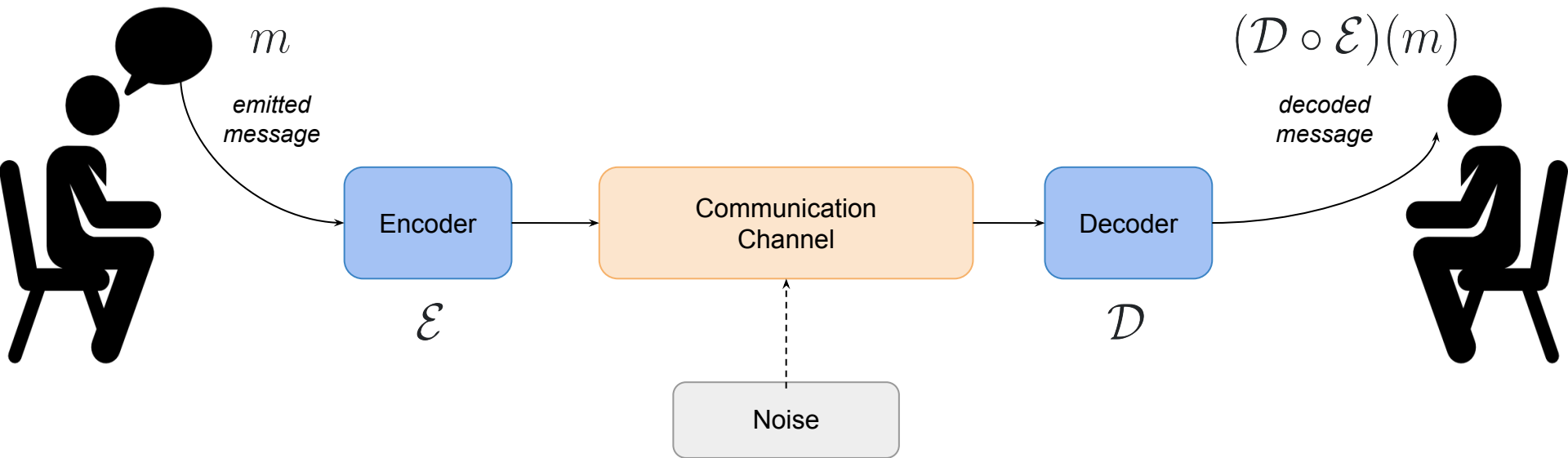
. Human communication



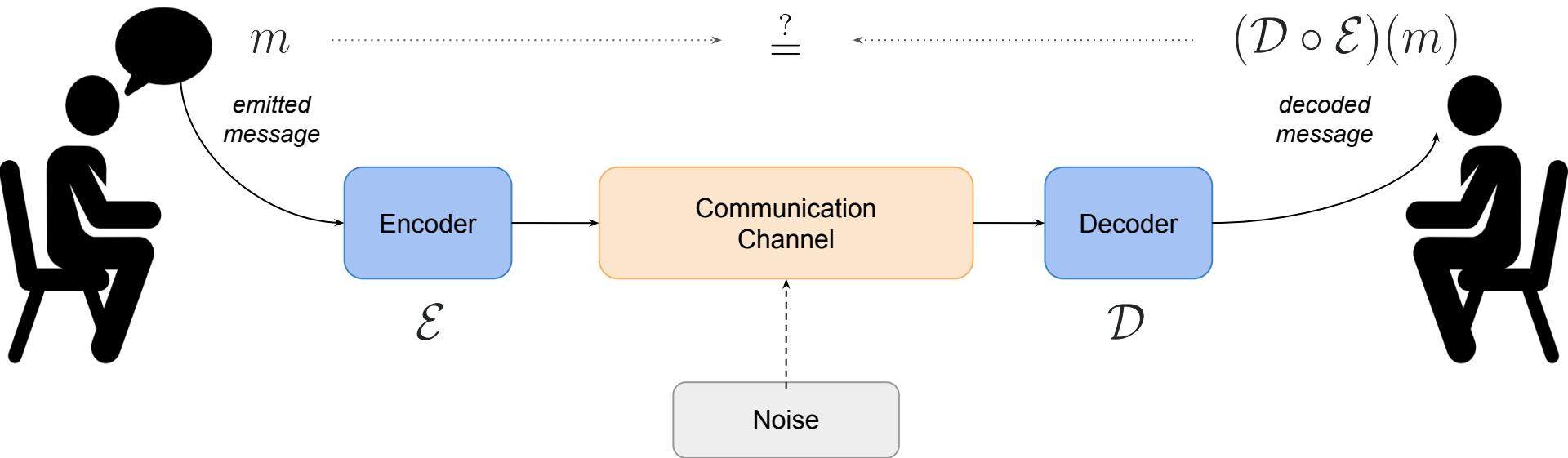
. Human communication



Human communication



Human communication



. Context – Temporality



. Context – Temporality



Context – Temporality

Responding

Listening

Following



Adapting

Coordinating

Synchronizing

Context – Temporality

Responding

Listening

Following



Adapting

Coordinating

Synchronizing

Aligning

Motivations

Uncovering the *temporal relationships* in audio and voice signals is of utmost importance for their *analysis*

When is a note played?

When is a word pronounced?

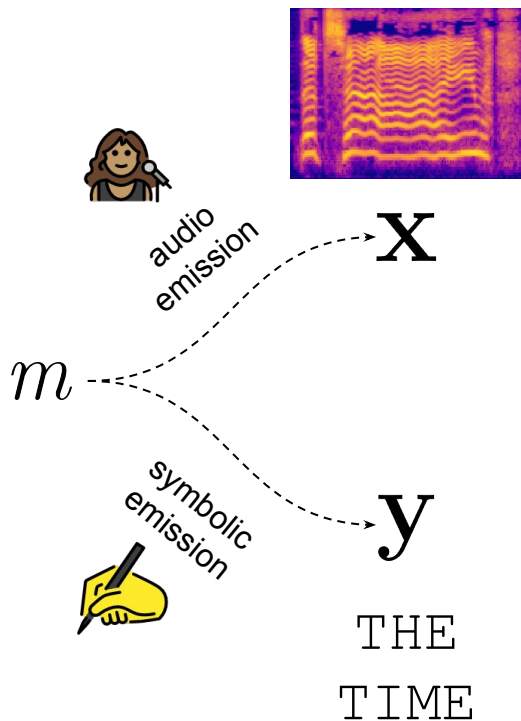
When is there a structural change in the audio?

.
. .
.

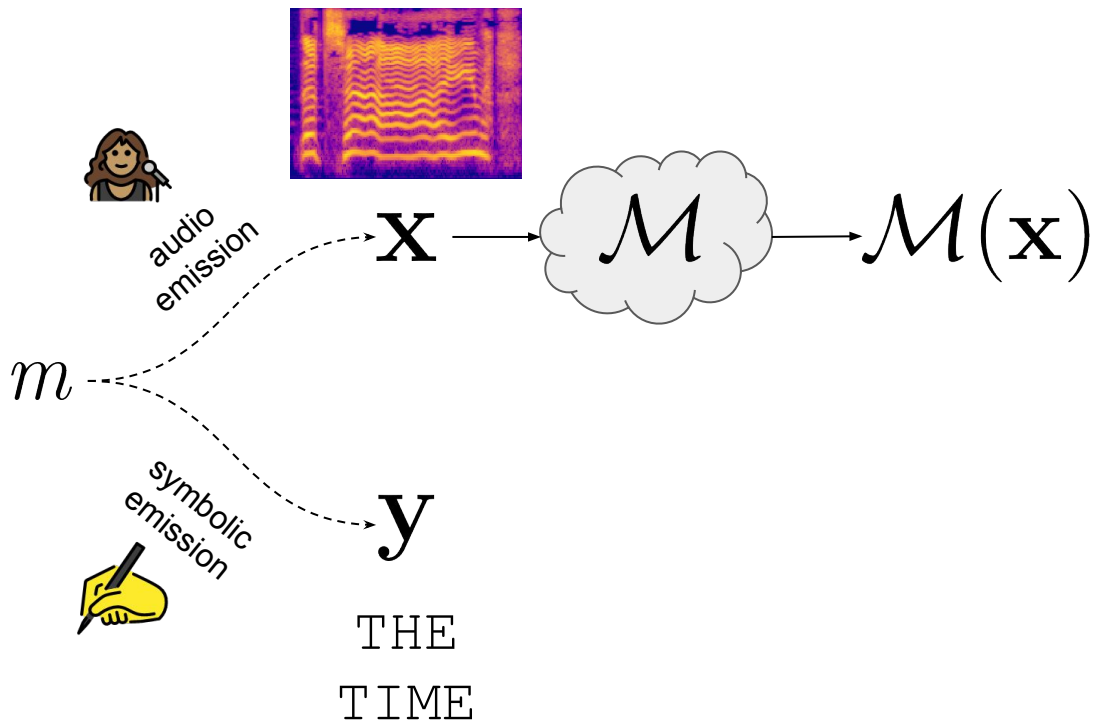
. Voice alignment

m

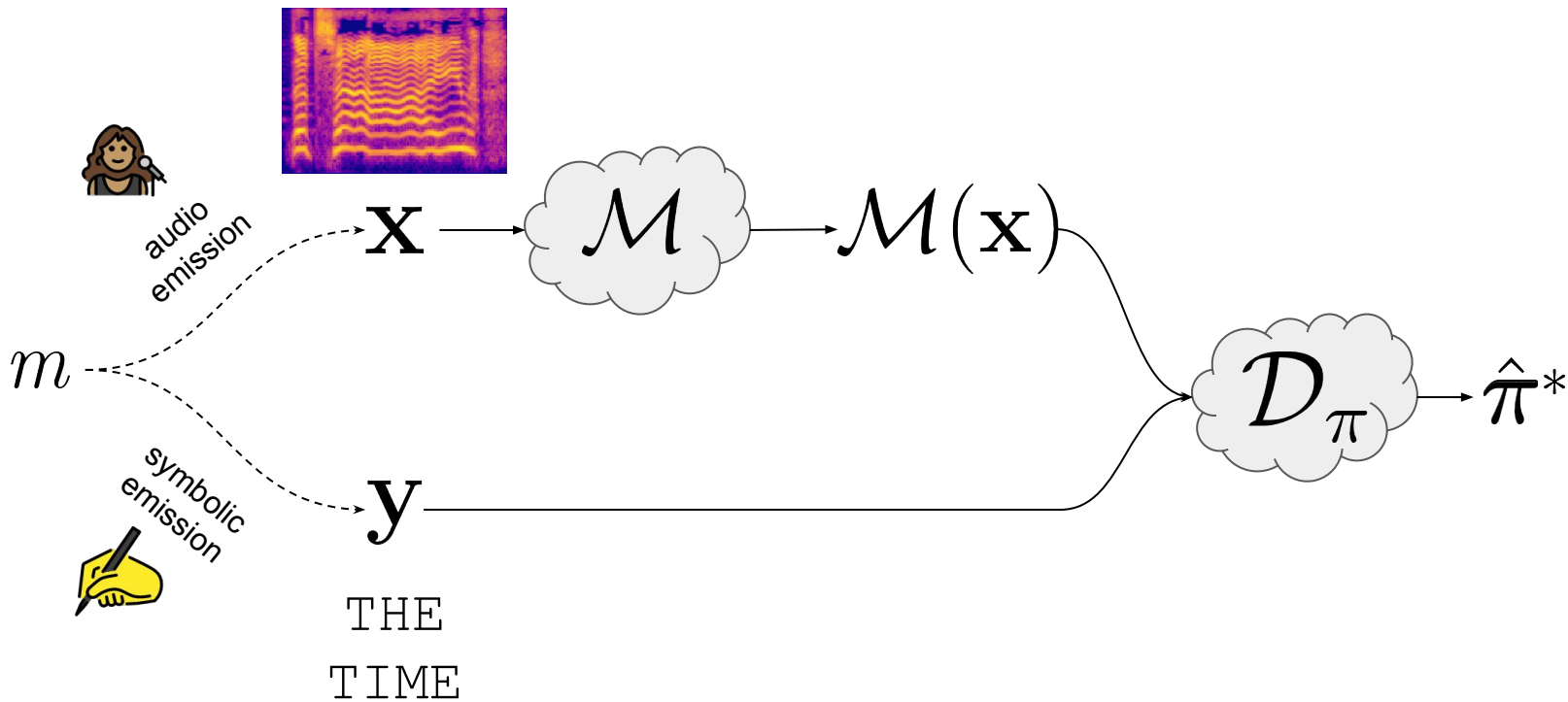
. Voice alignment



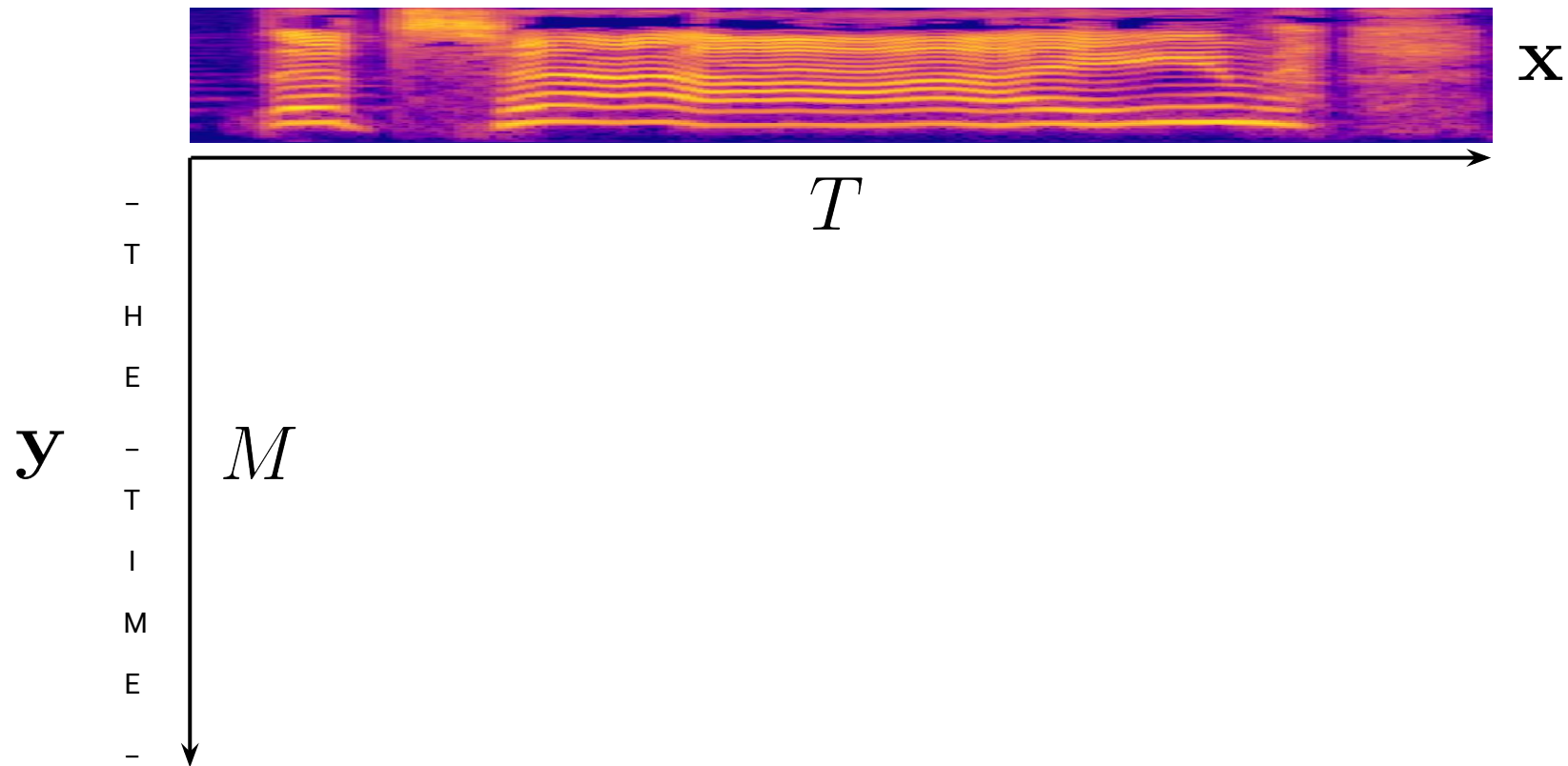
. Voice alignment



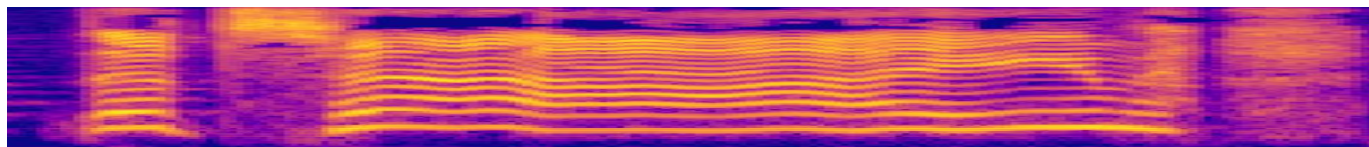
. Voice alignment



. Optimal pathway



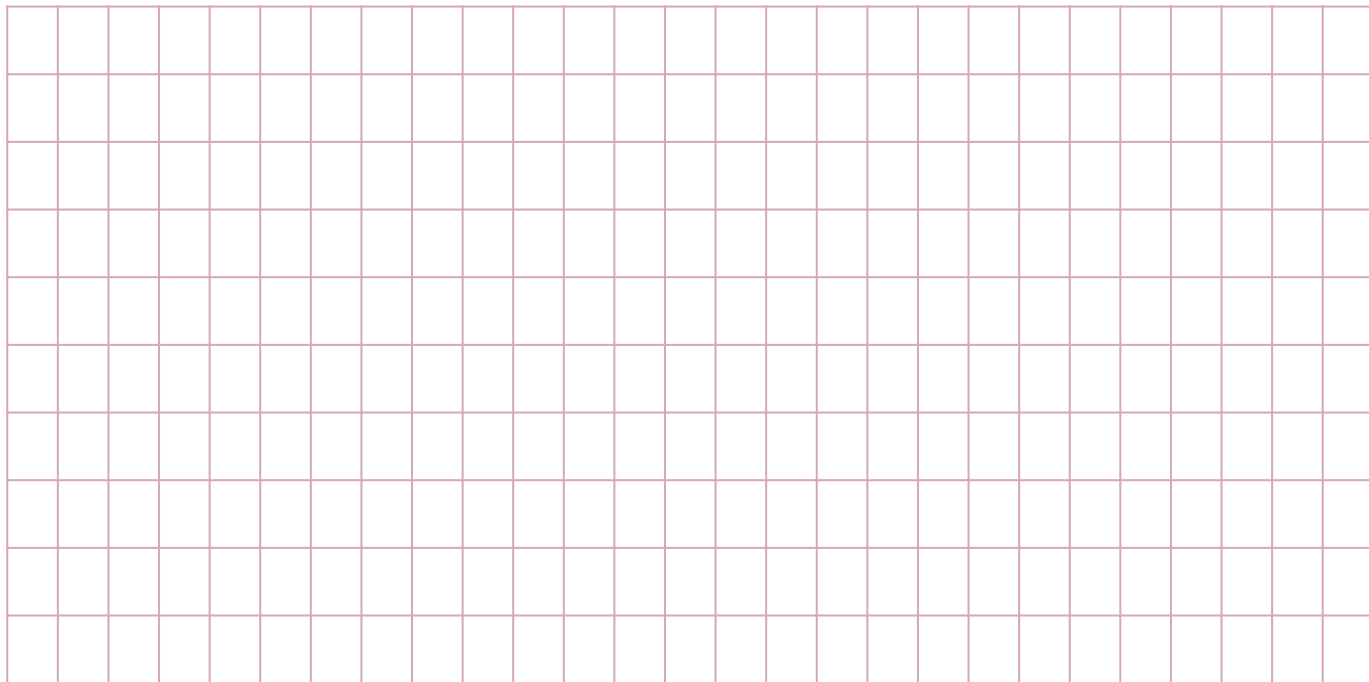
. Optimal pathway



x

y

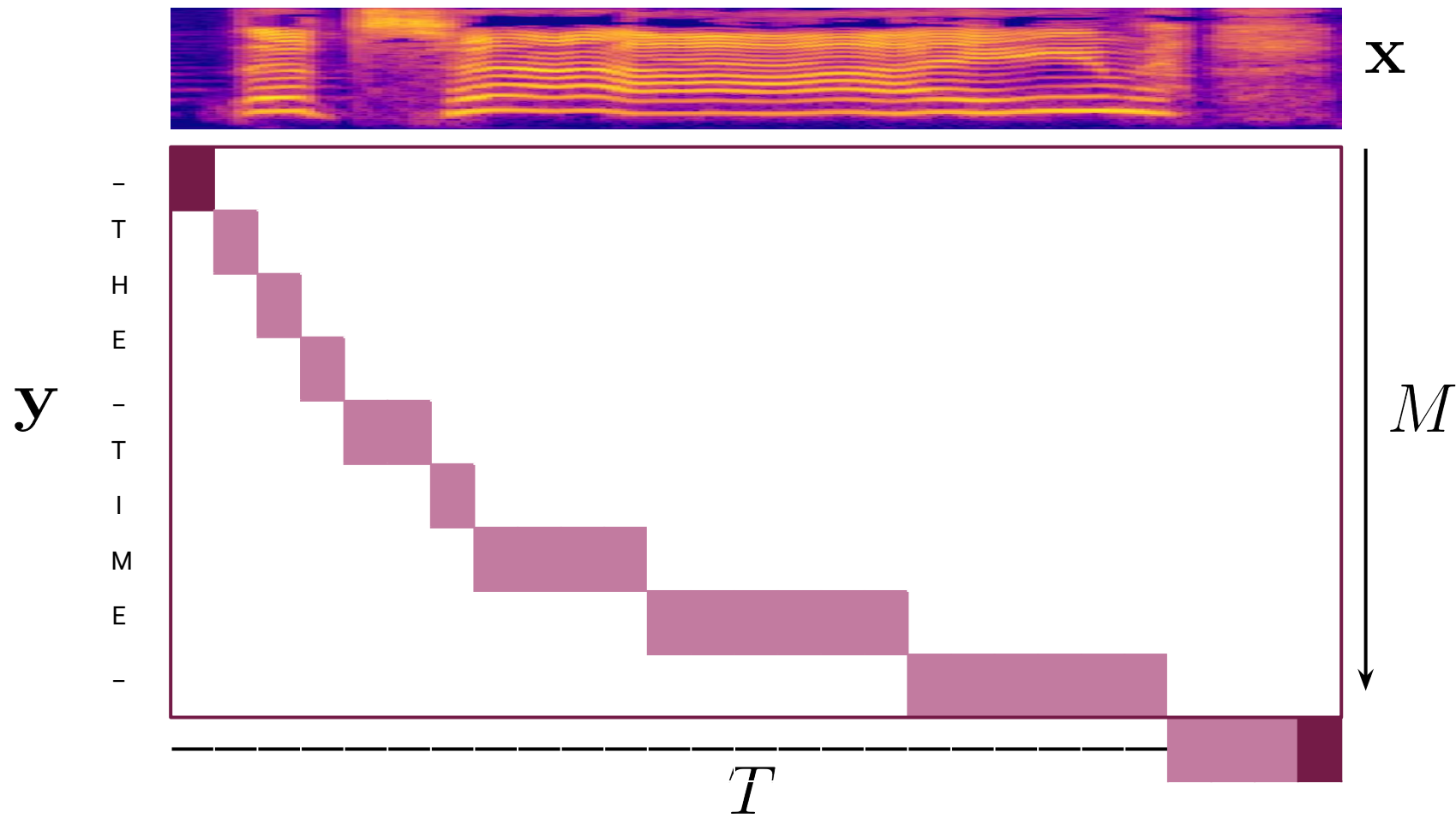
-
T
H
E
-
T
I
M
E
-



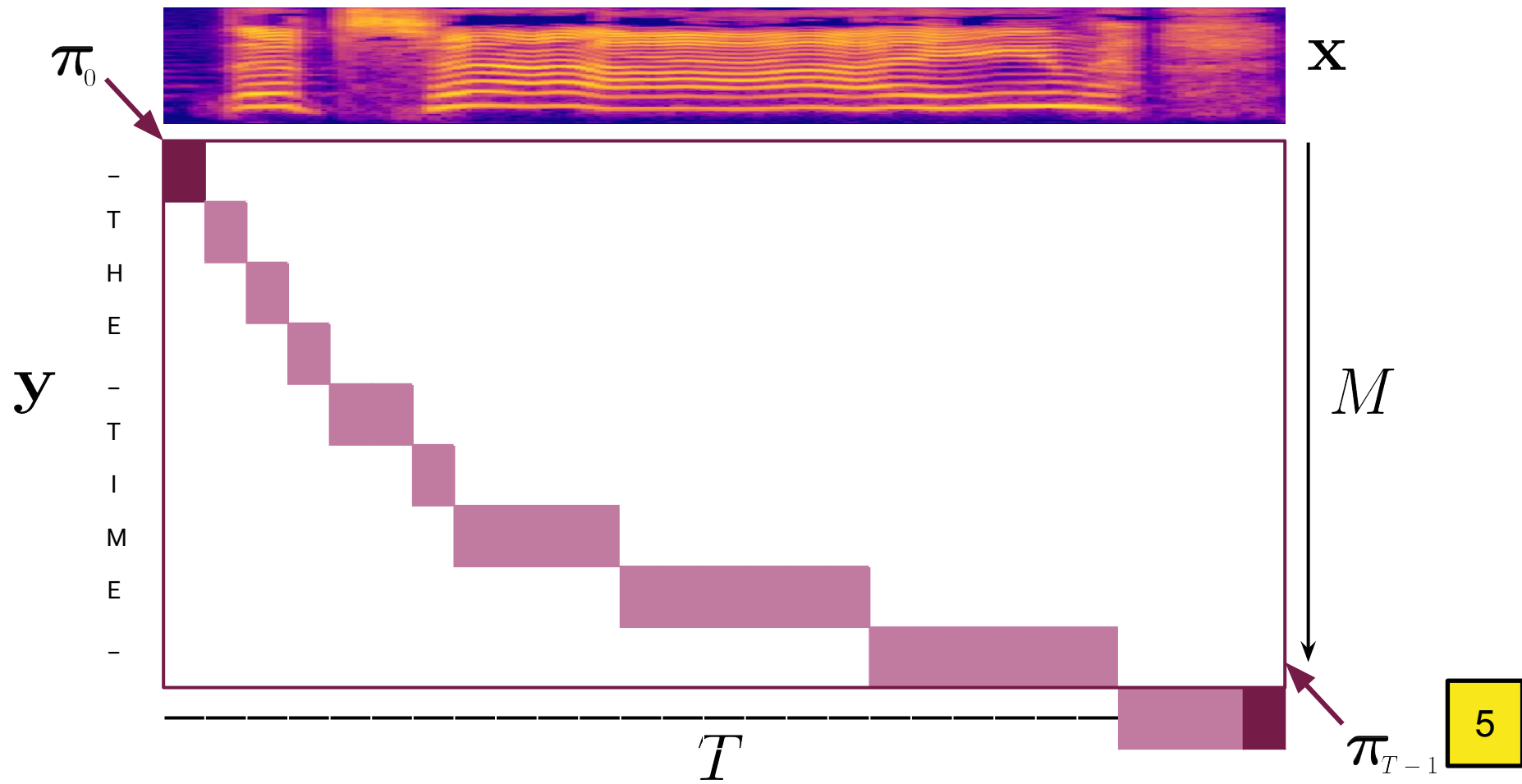
M

T

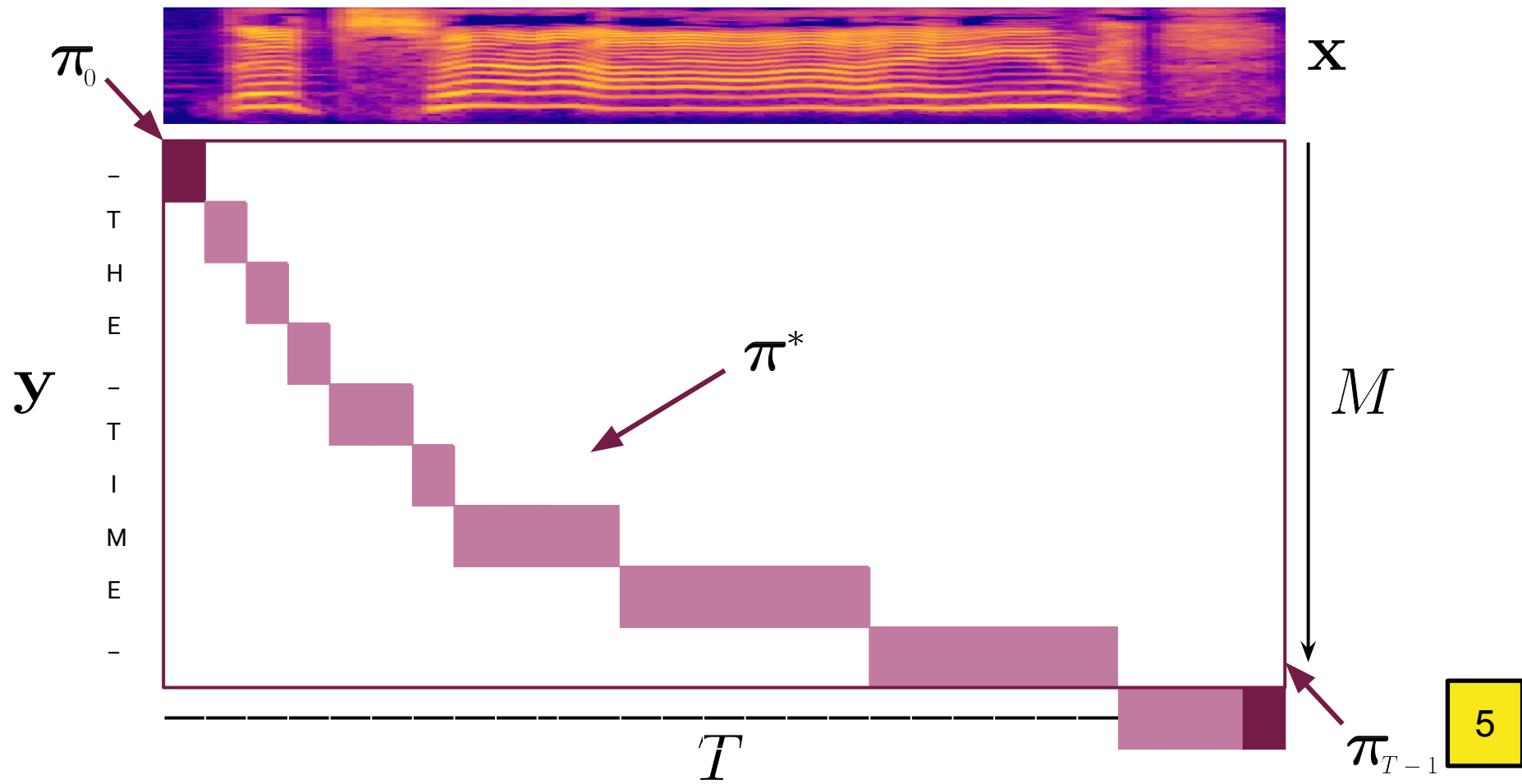
Optimal pathway



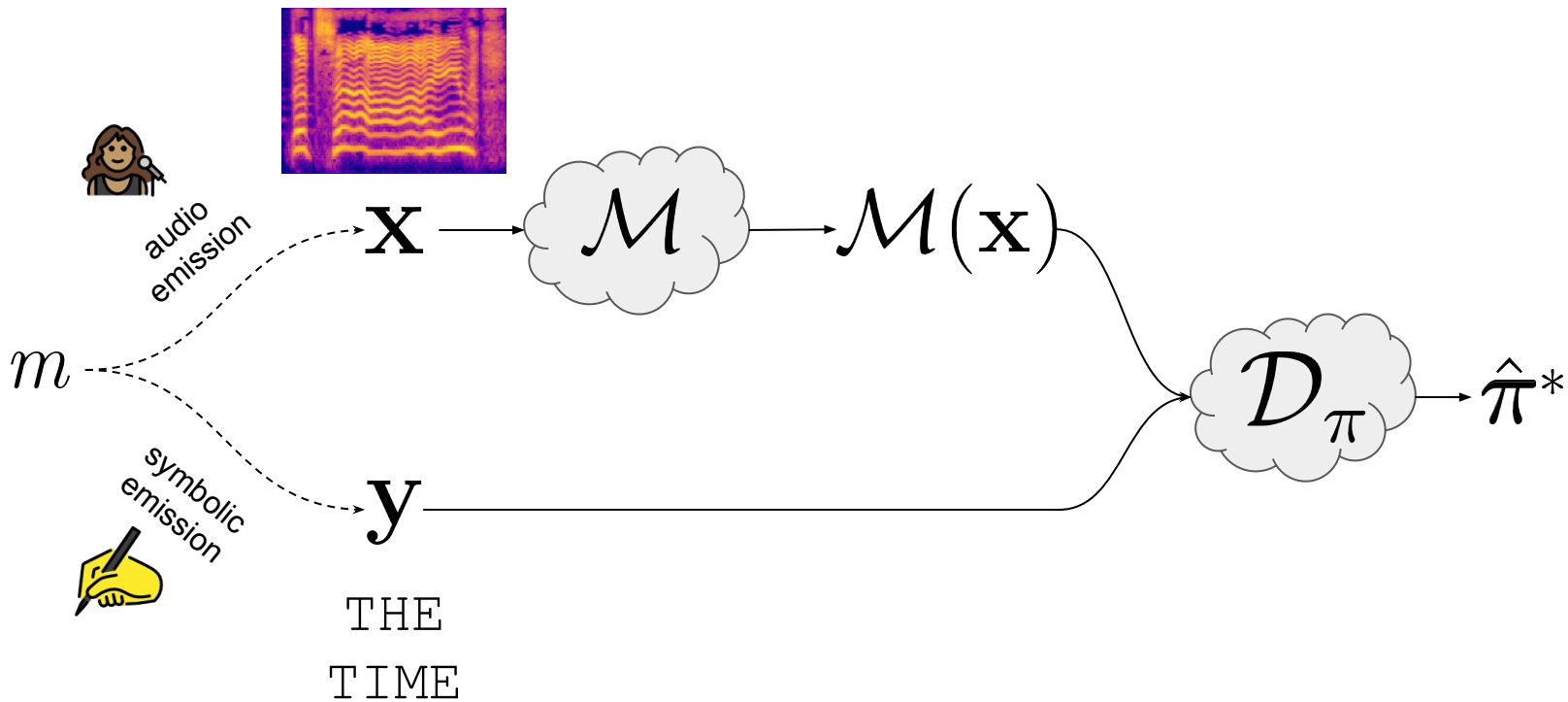
Optimal pathway



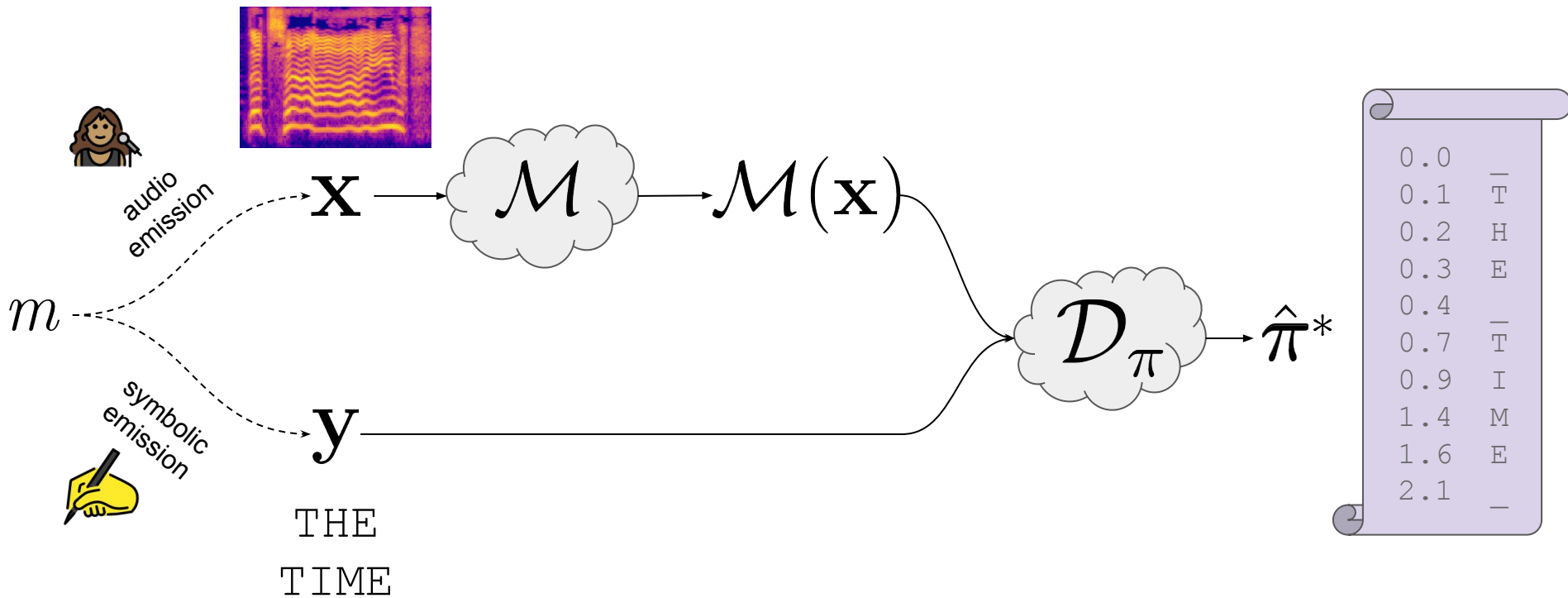
Optimal pathway



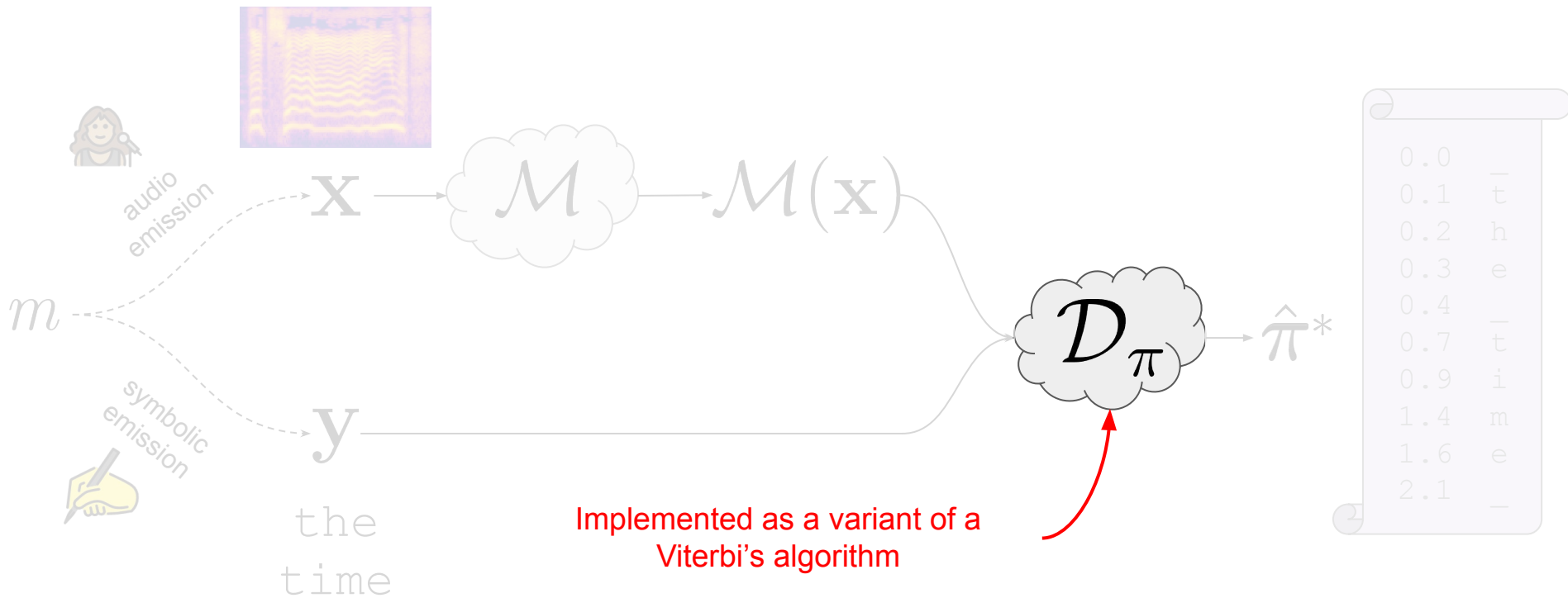
. Overview



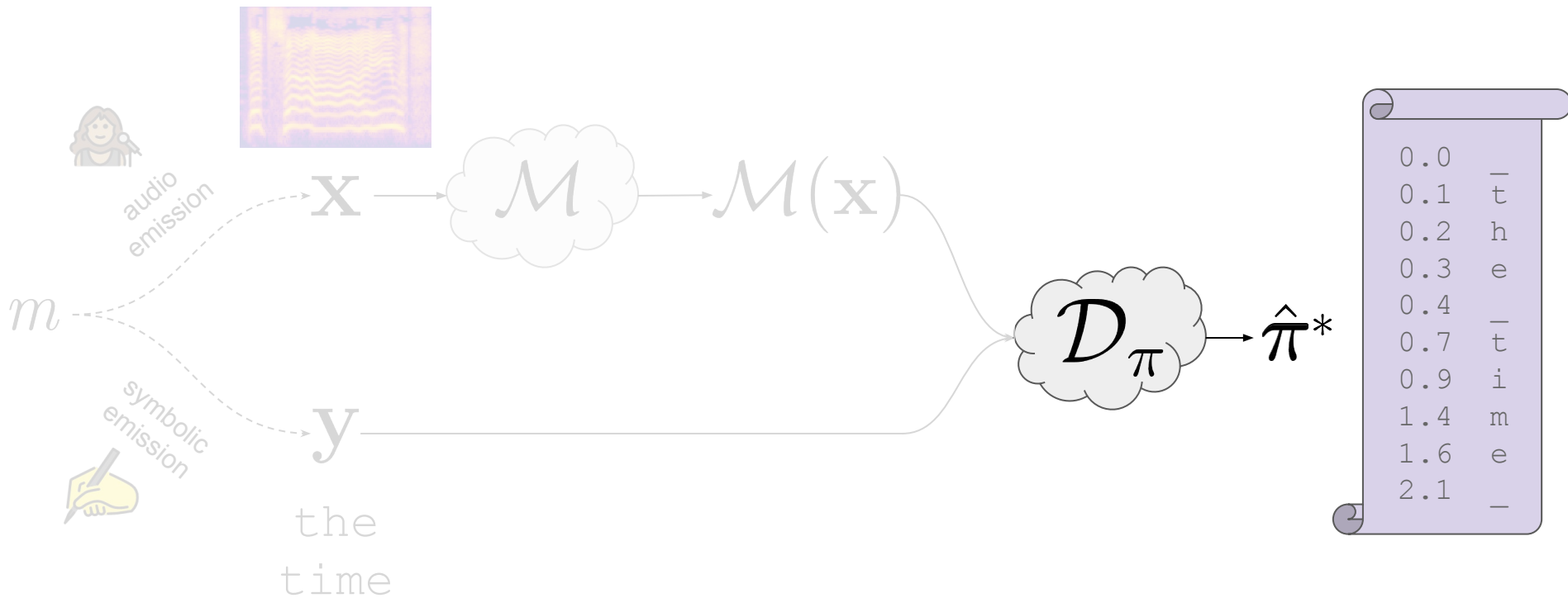
Overview



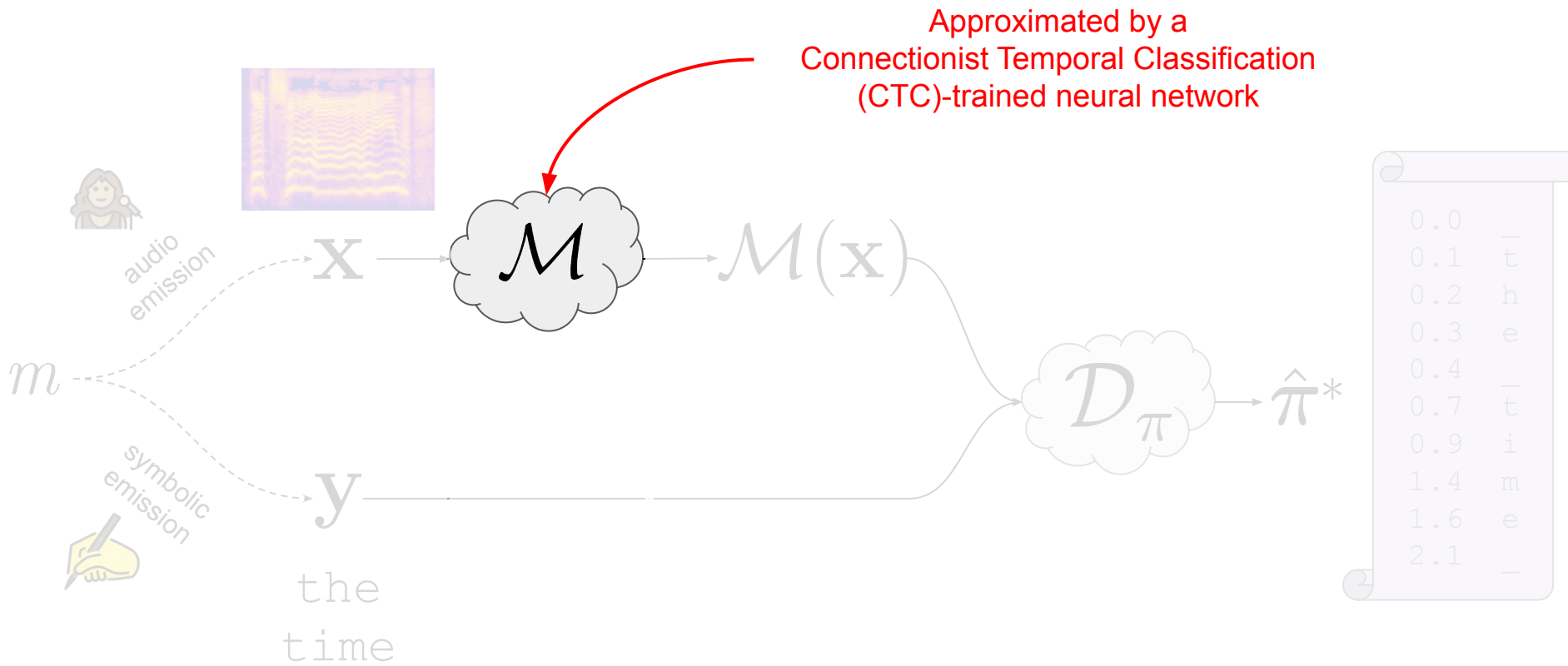
. Path decoding



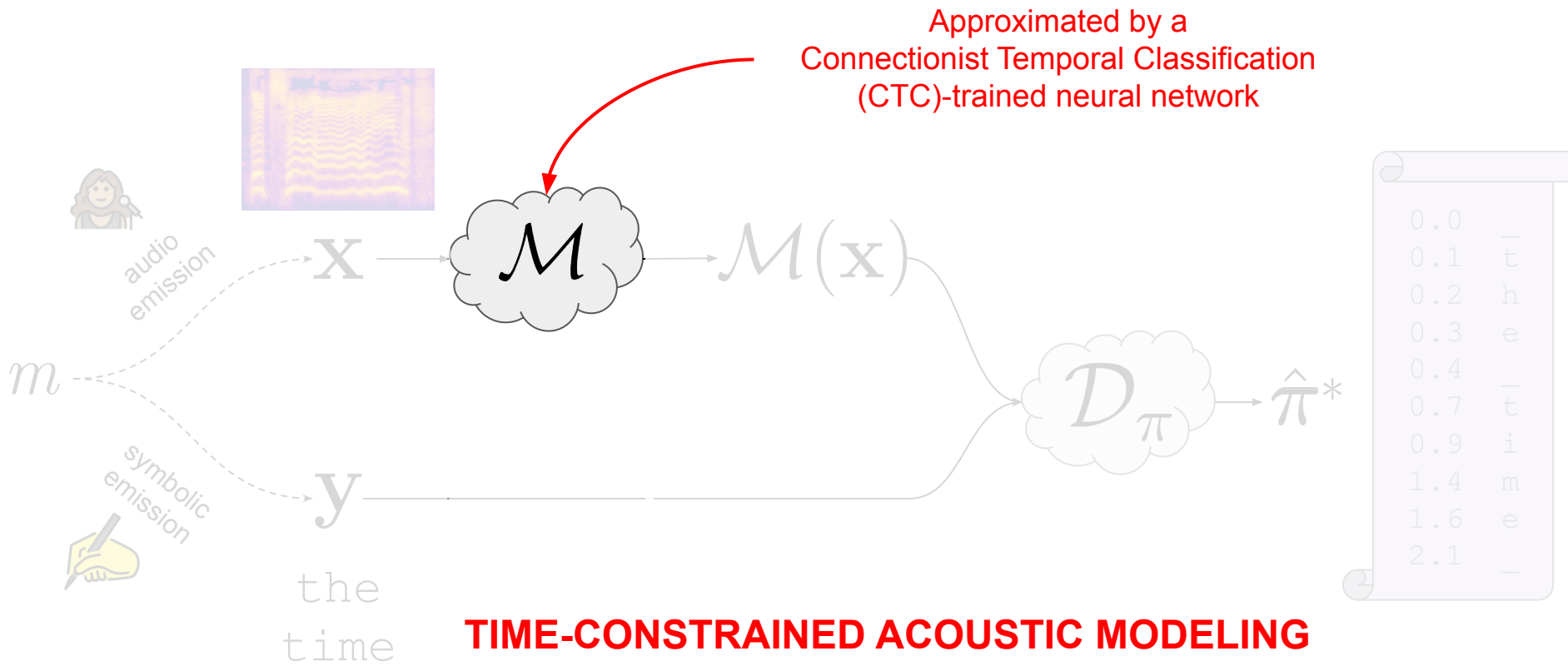
. Path decoding



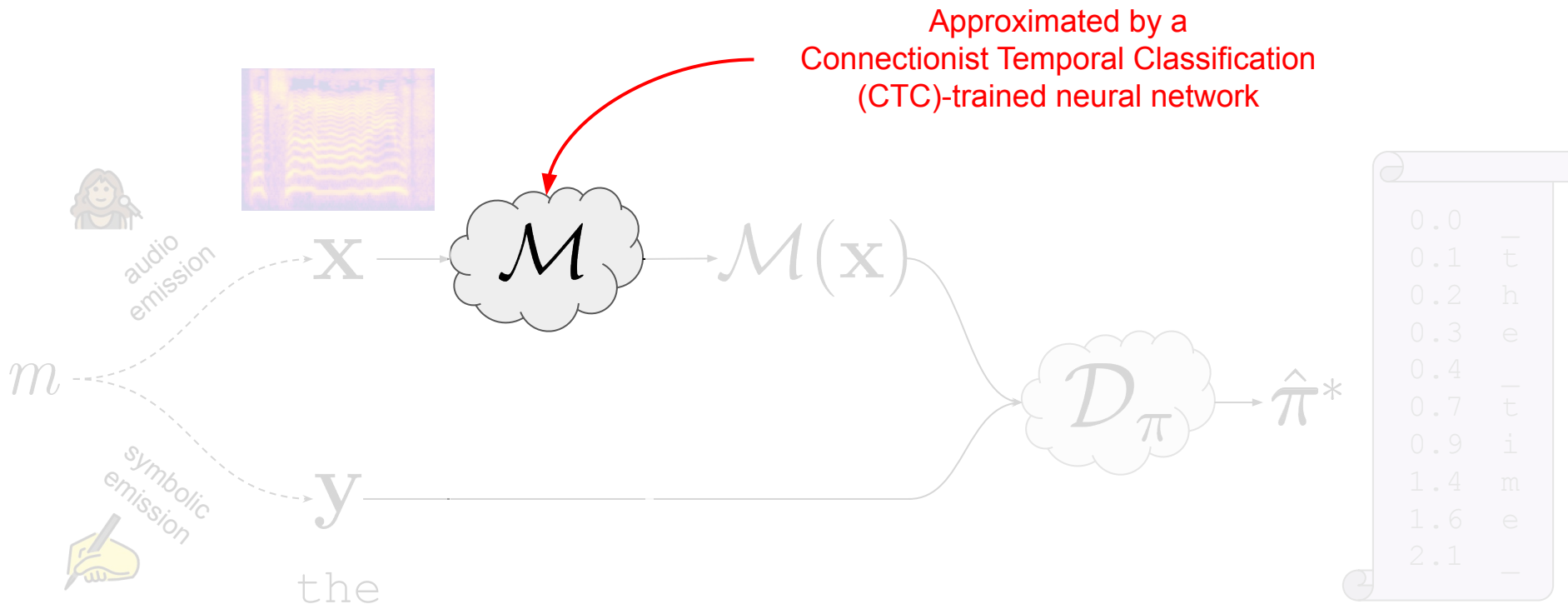
. Acoustic modeling



. Acoustic modeling



. Acoustic modeling



**TIME-CONSTRAINED ACOUSTIC MODELING
DETAILS ON POSTER :)**