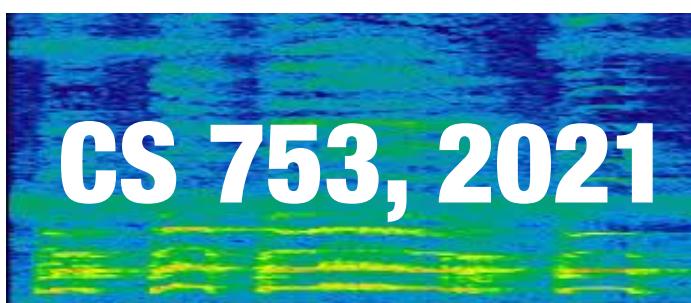


WFSTs in ASR

Lecture 4b



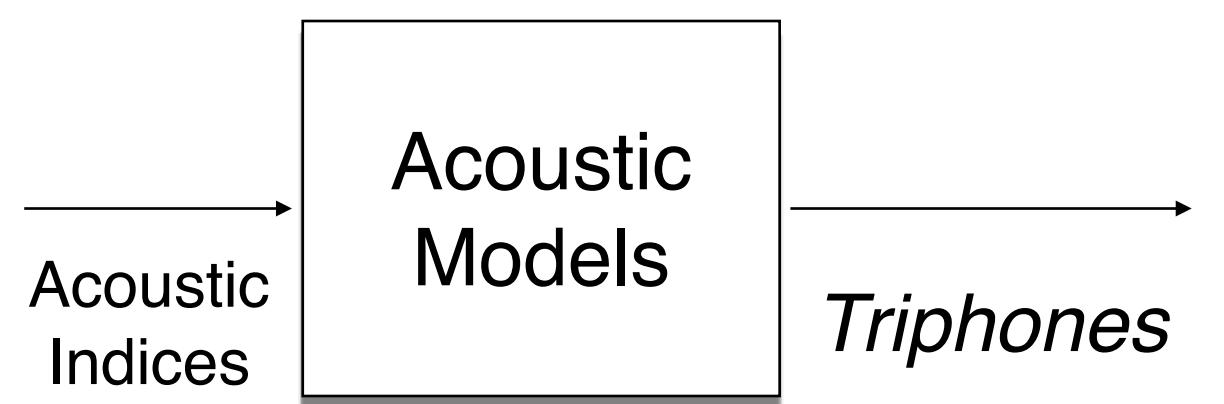
Instructor: Preethi Jyothi, IITB

WFST-based ASR System

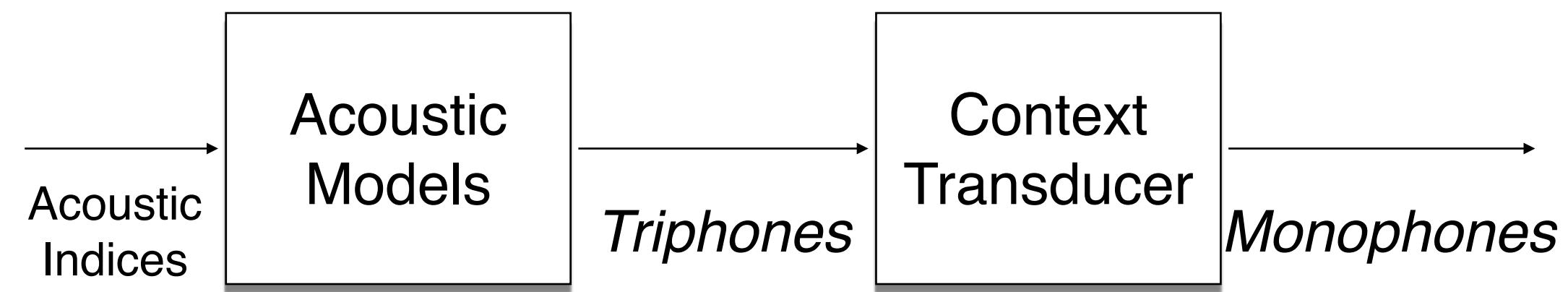
WFST-based ASR System

→
Acoustic
Indices

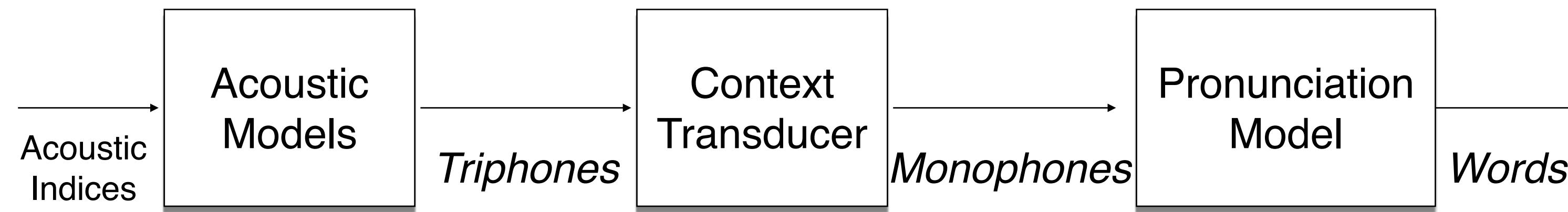
WFST-based ASR System



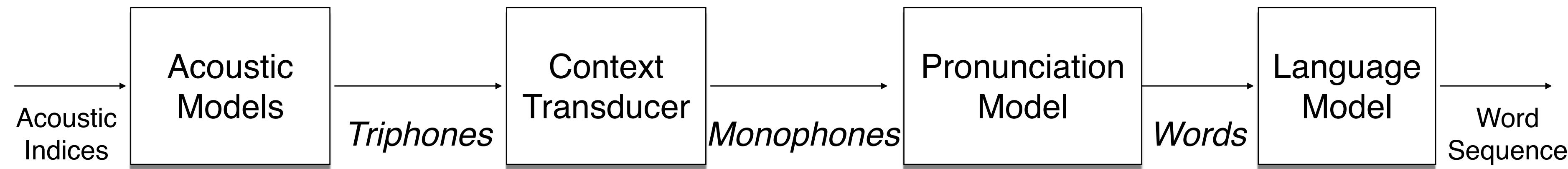
WFST-based ASR System



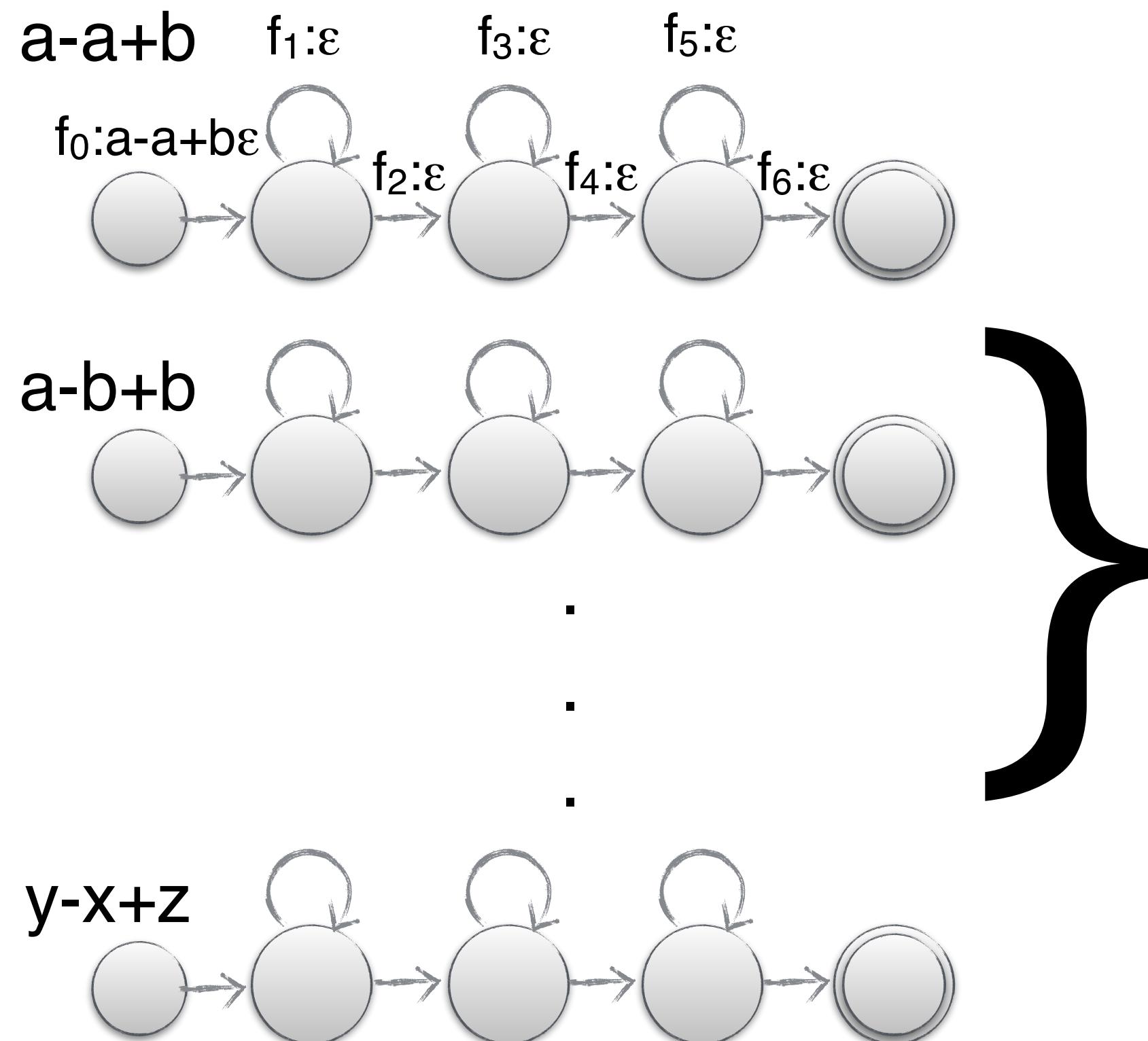
WFST-based ASR System



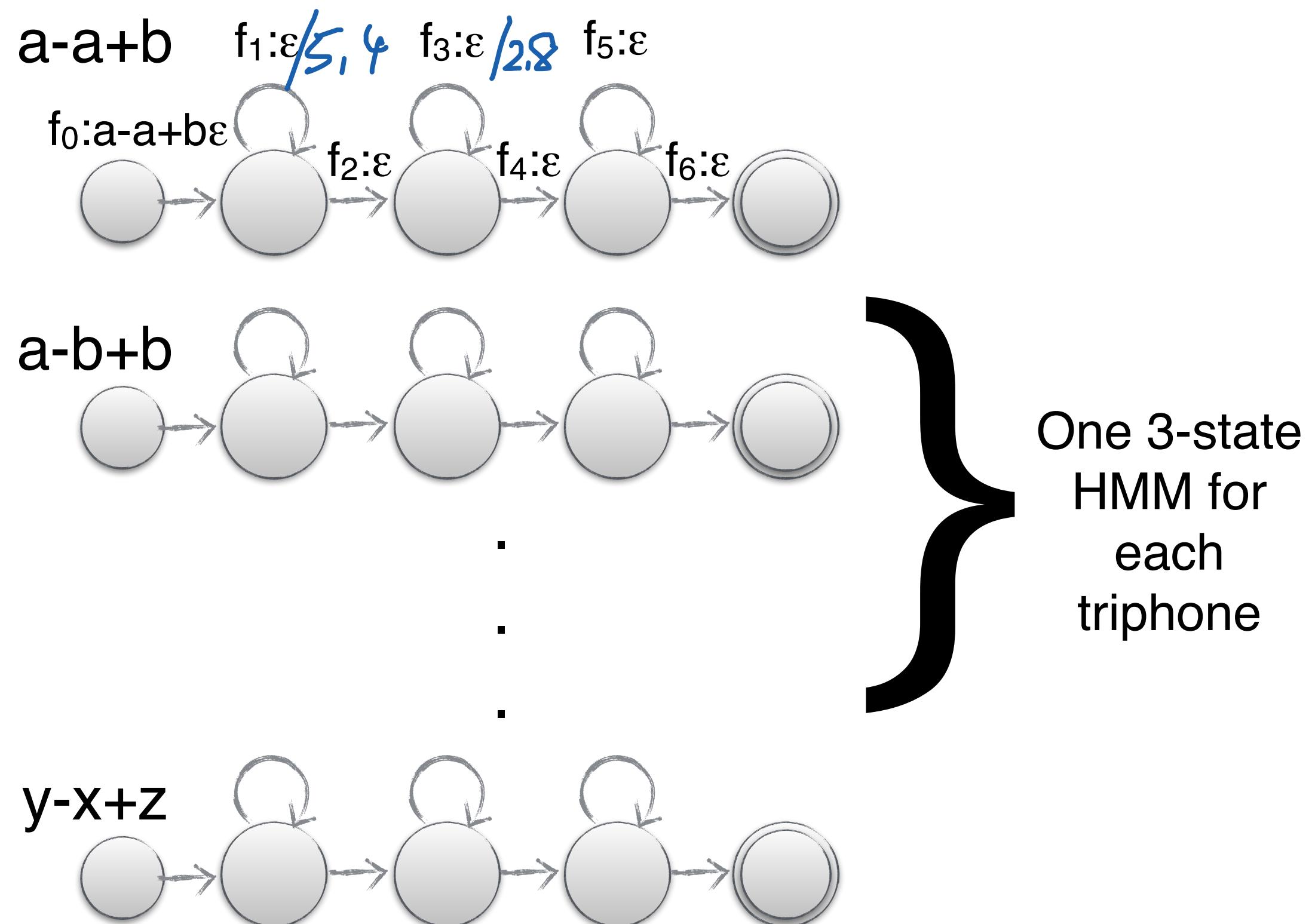
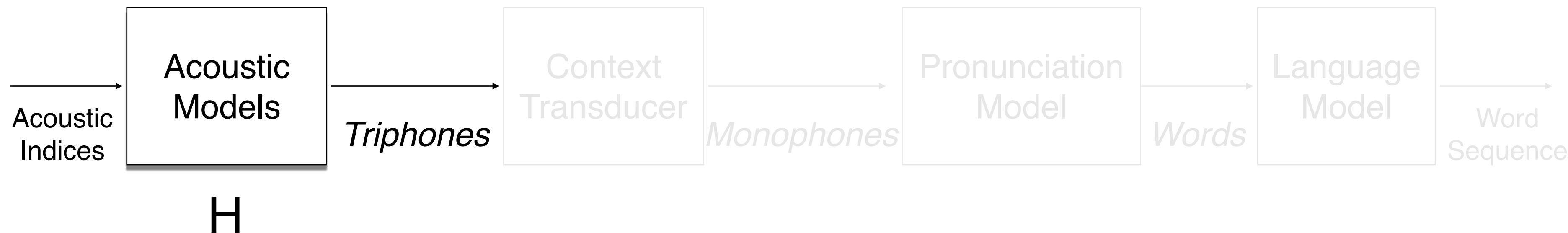
WFST-based ASR System



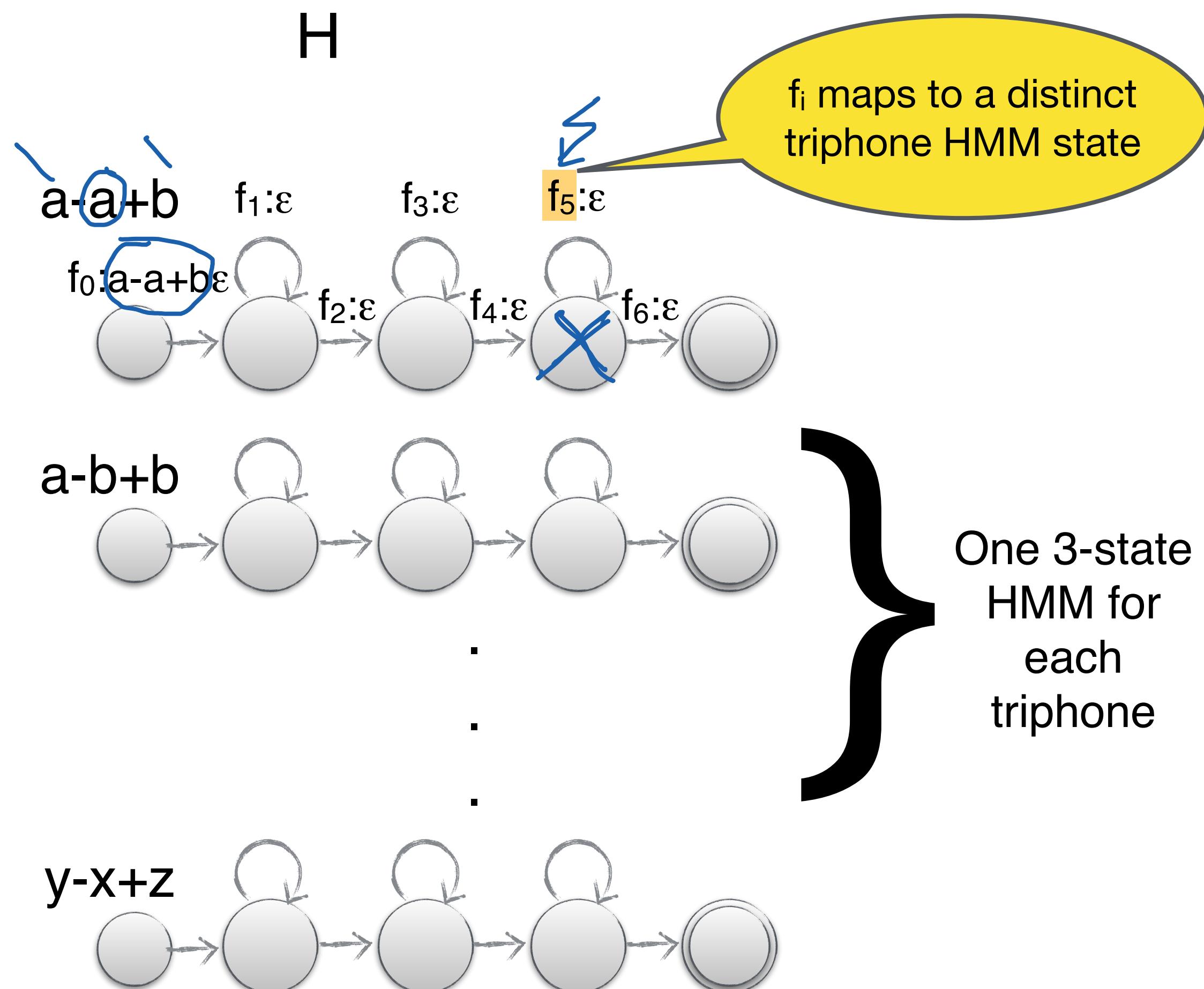
WFST-based ASR System



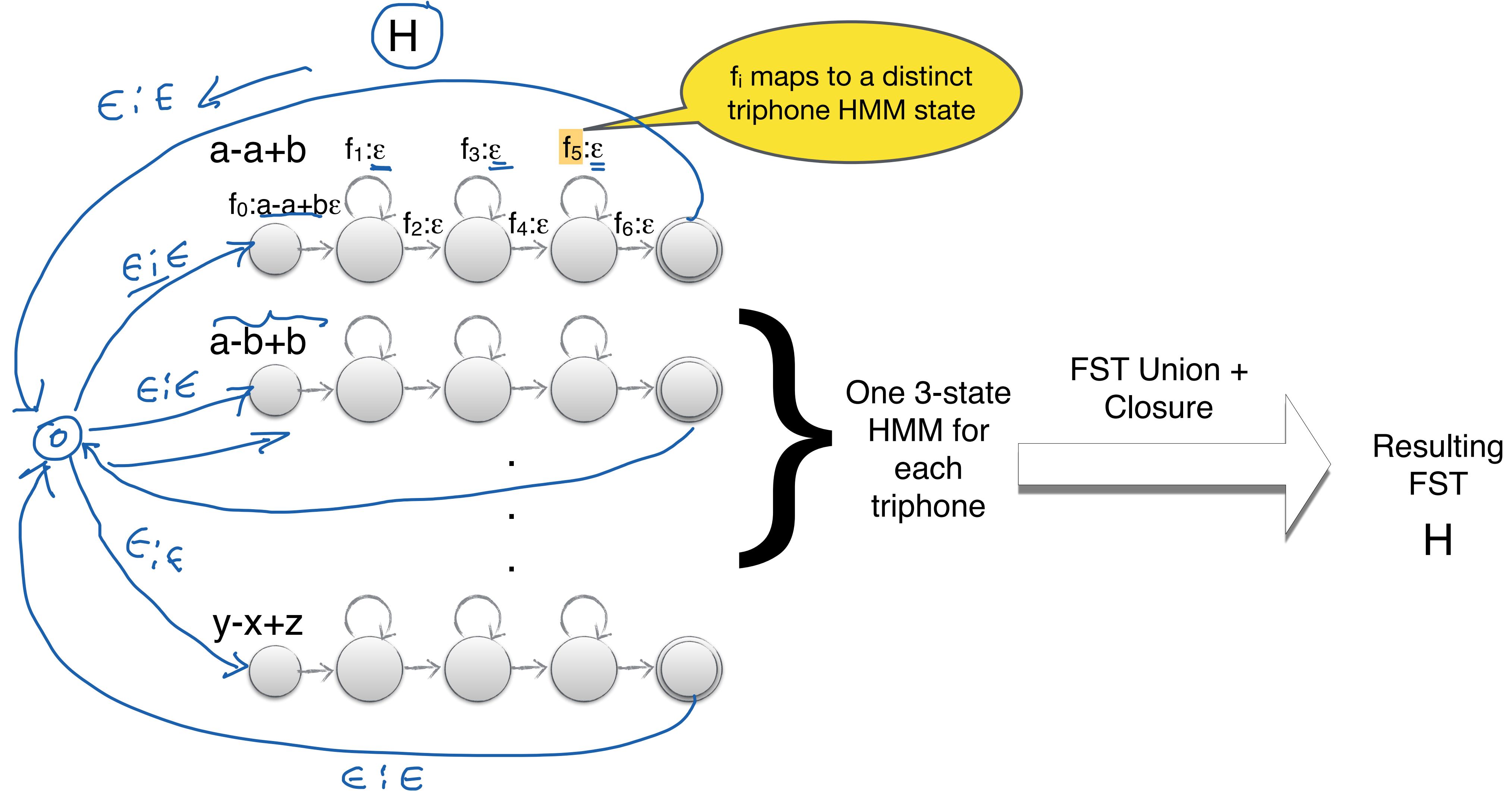
WFST-based ASR System



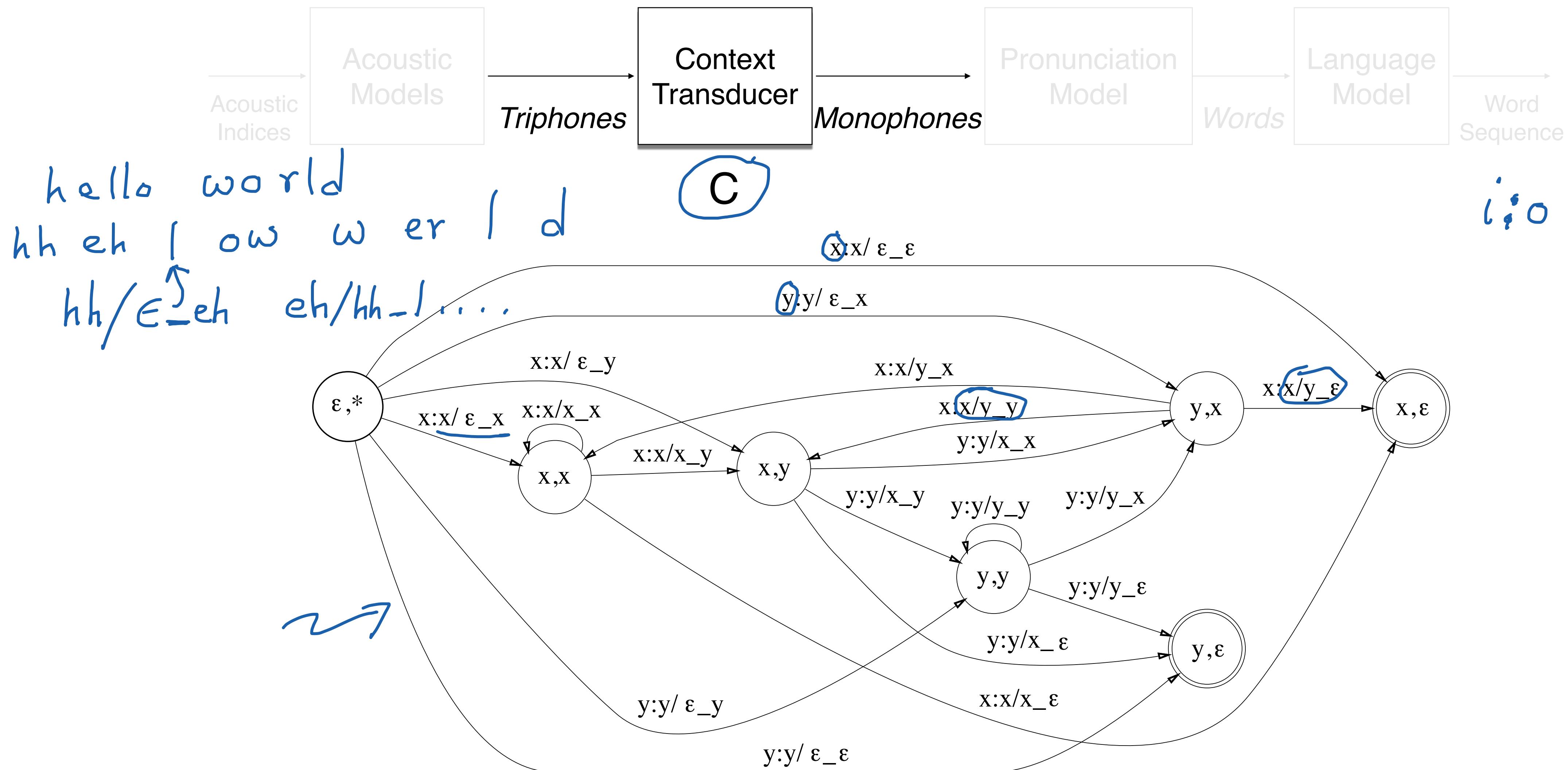
WFST-based ASR System



WFST-based ASR System

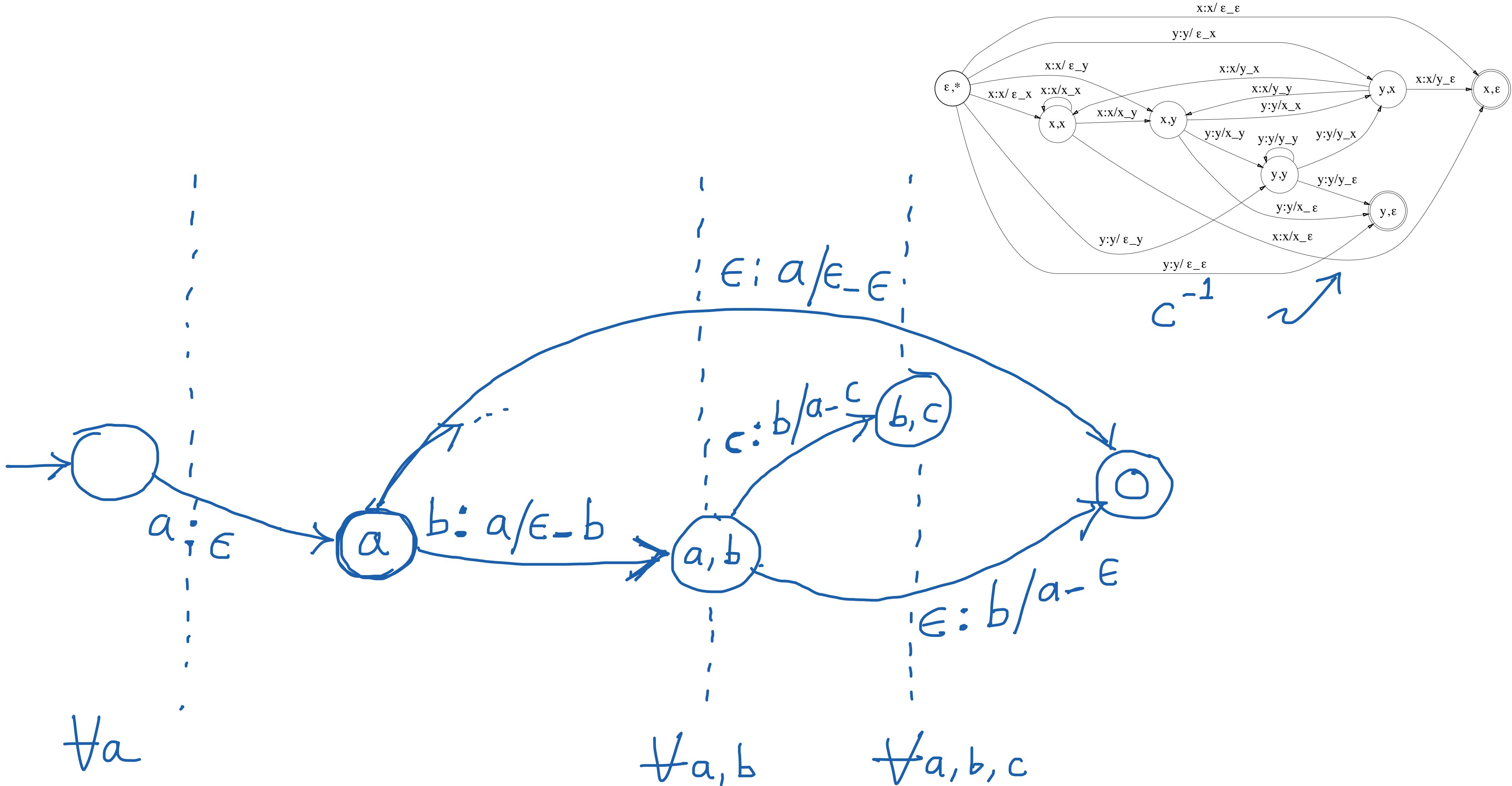


WFST-based ASR System

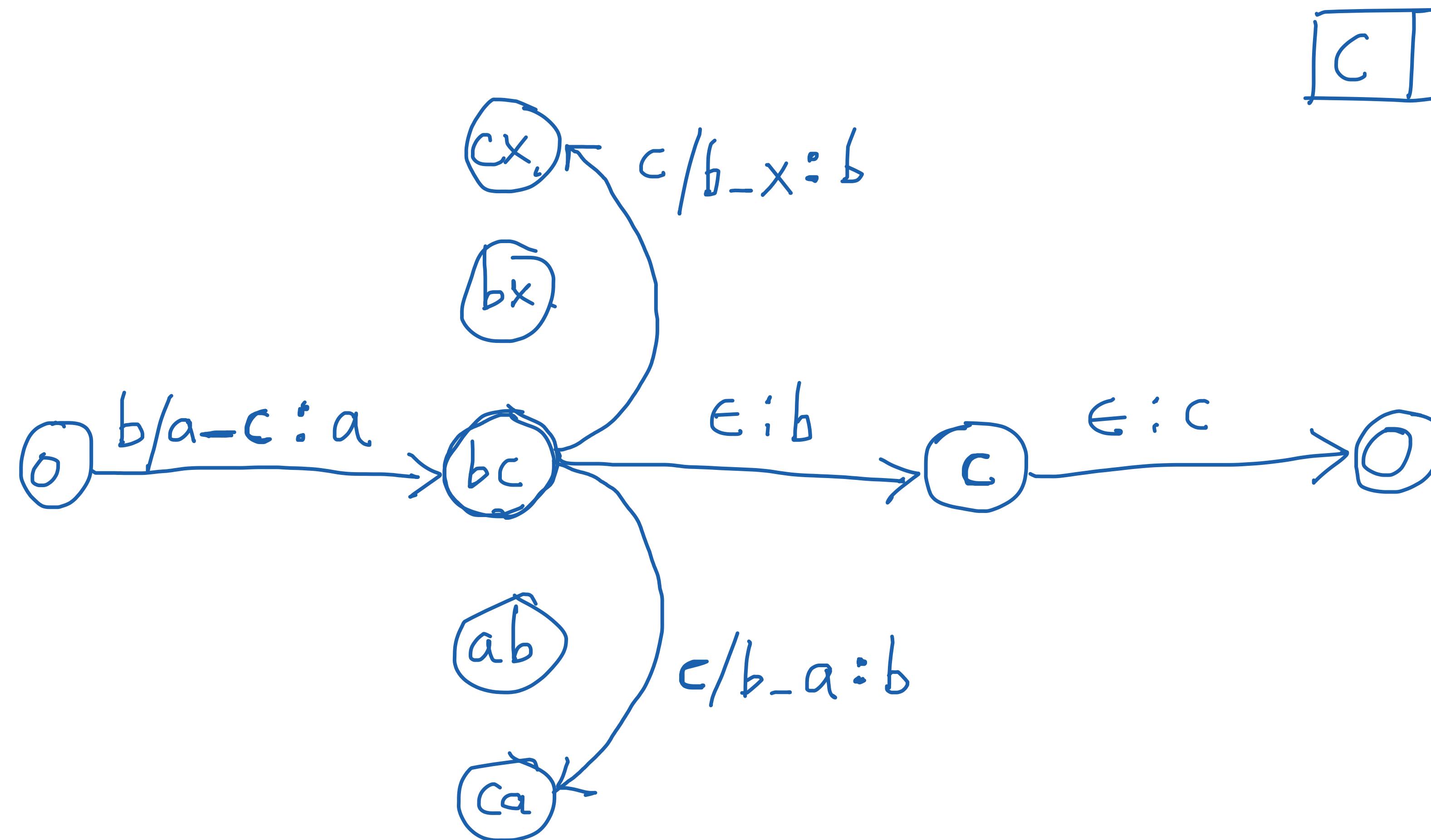


C⁻¹: Arc labels: “monophone : phone / left-context_right-context”

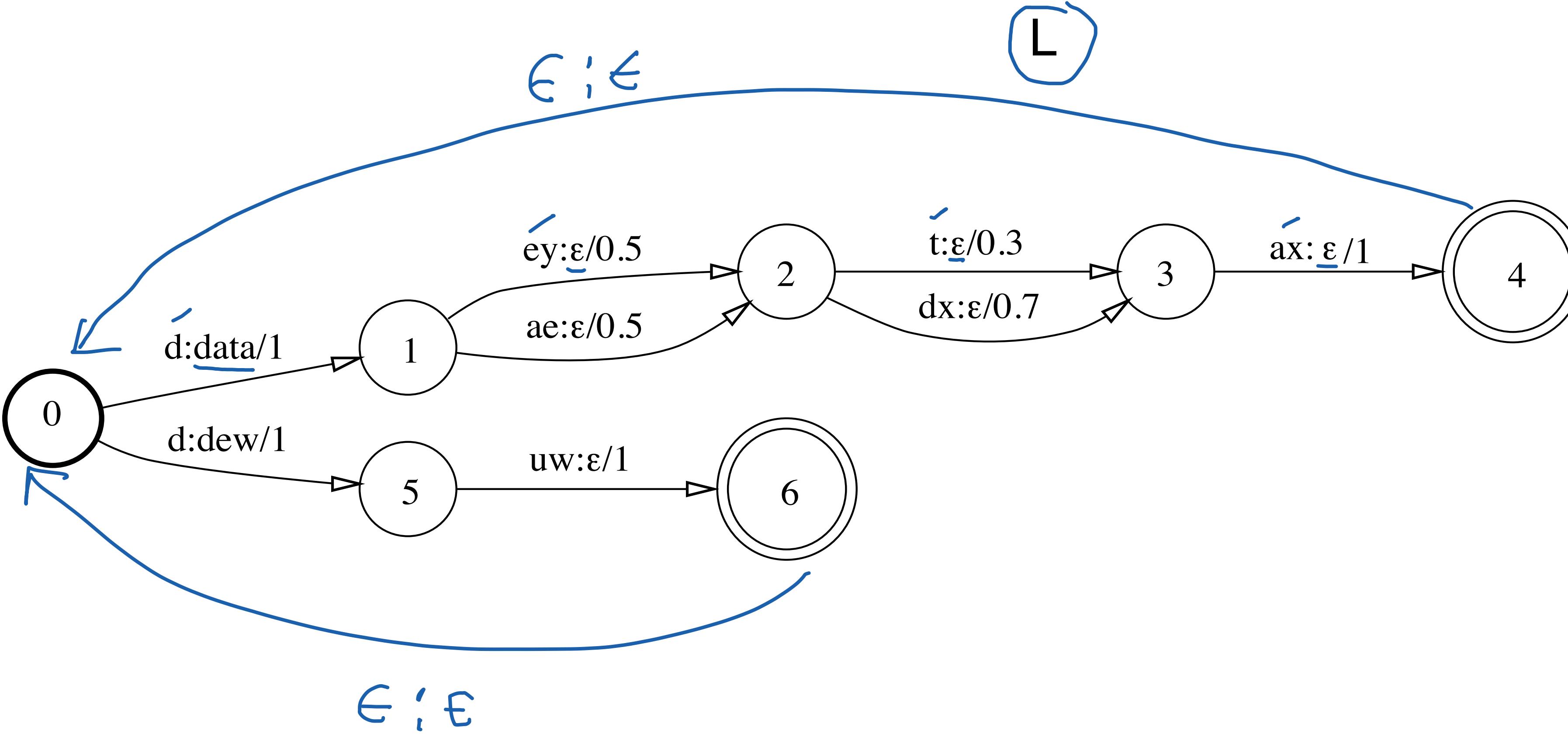
C: Context Dependency Transducer



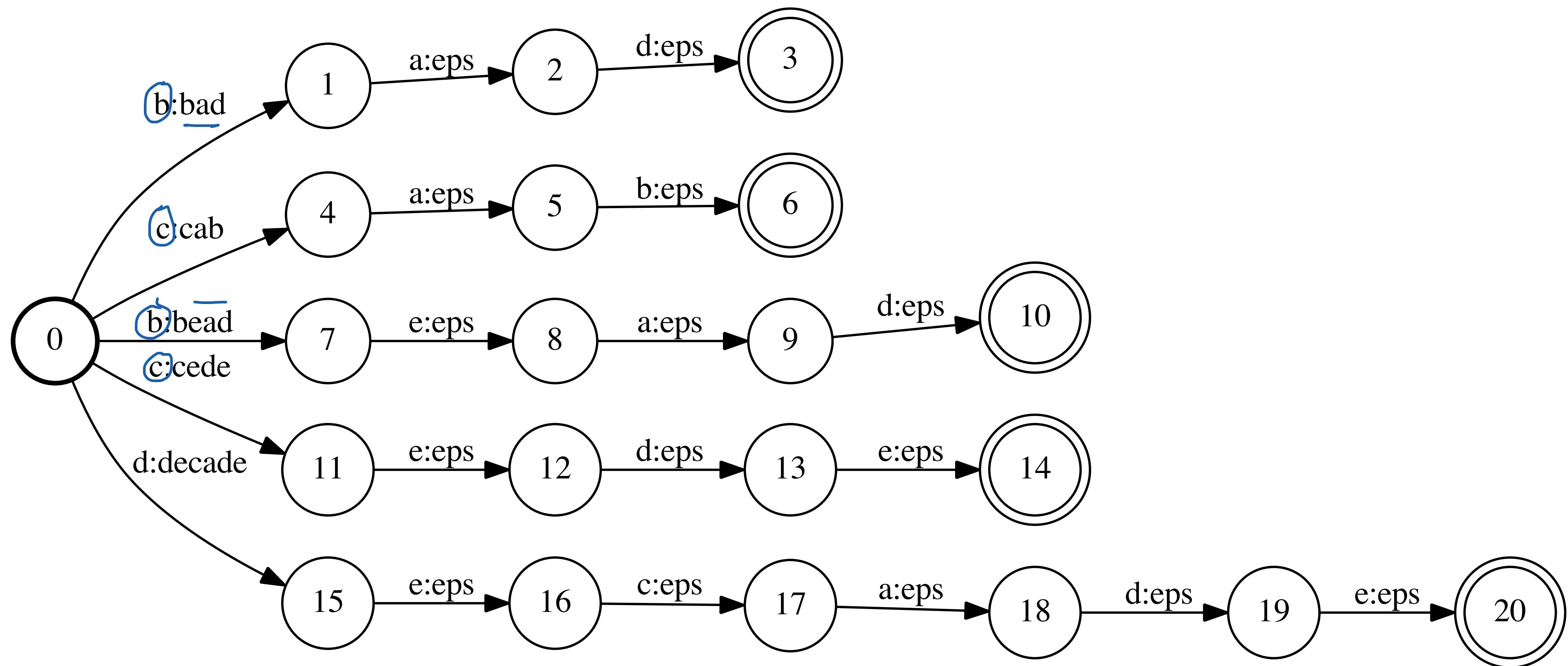
C: Context Dependency Transducer



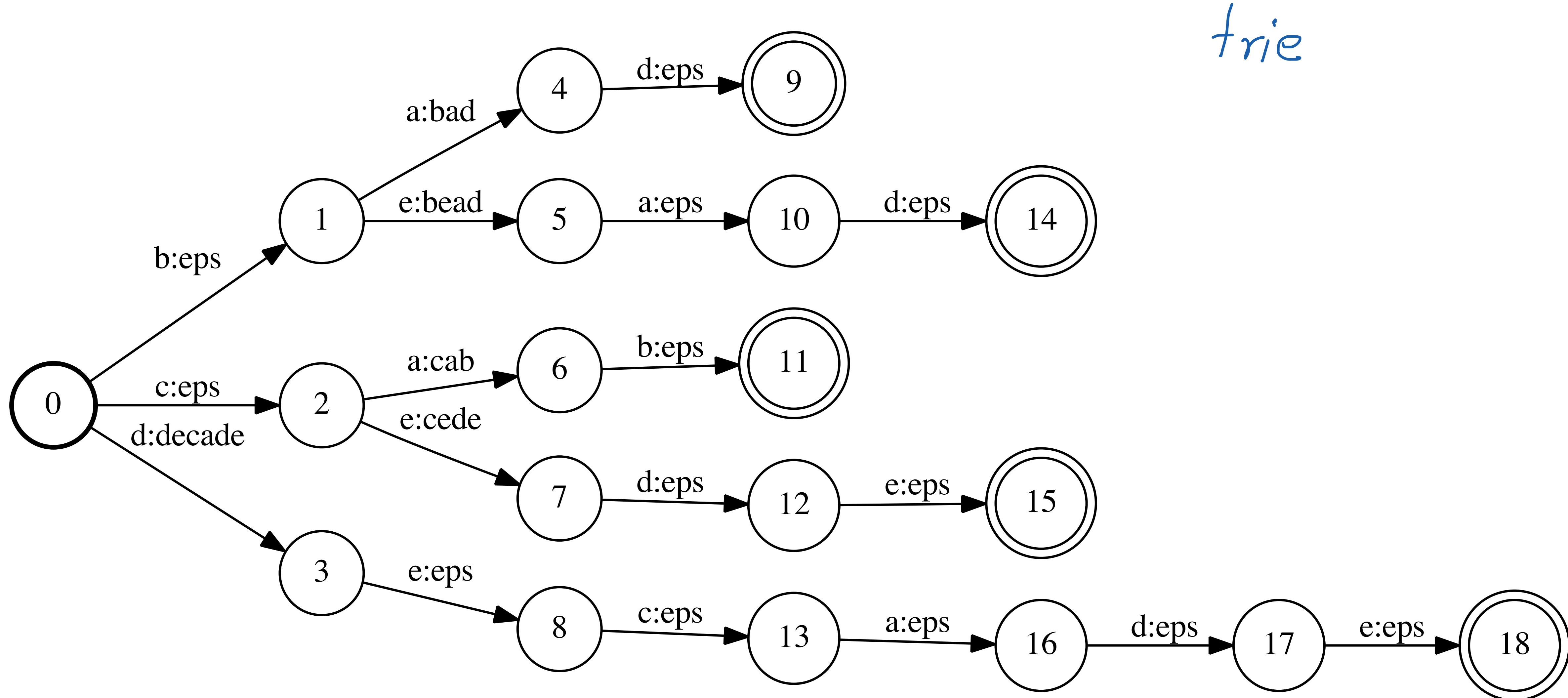
WFST-based ASR System



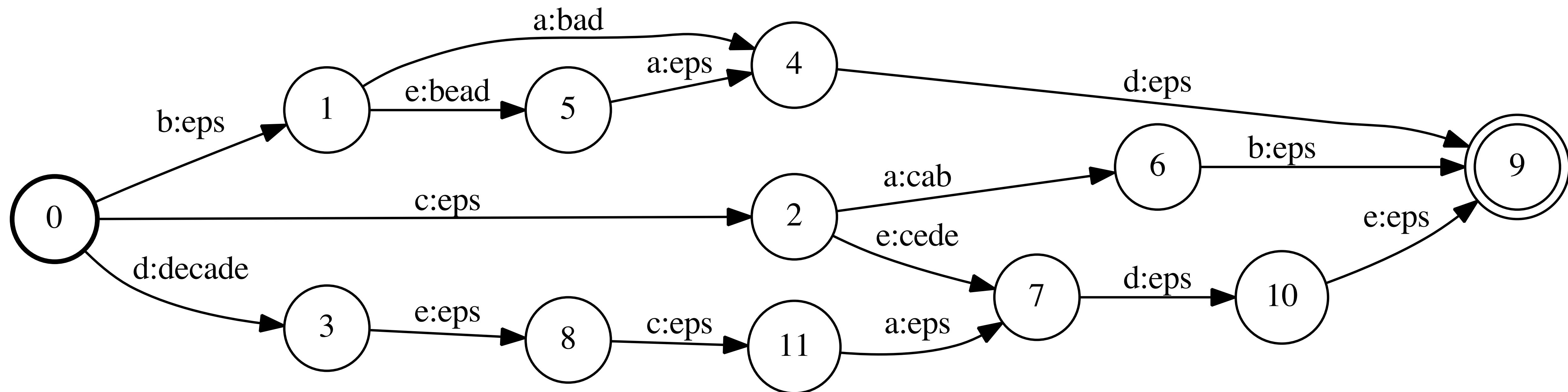
Example: Dictionary WFST



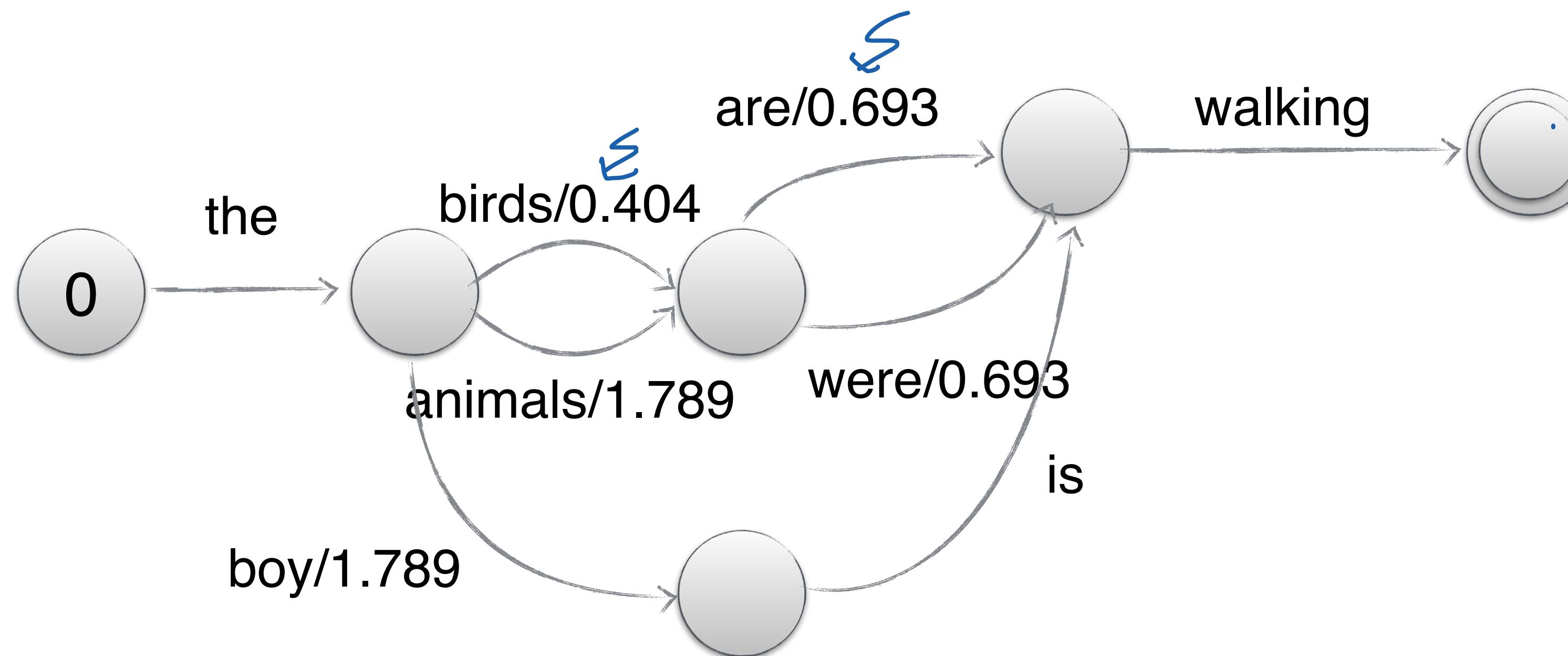
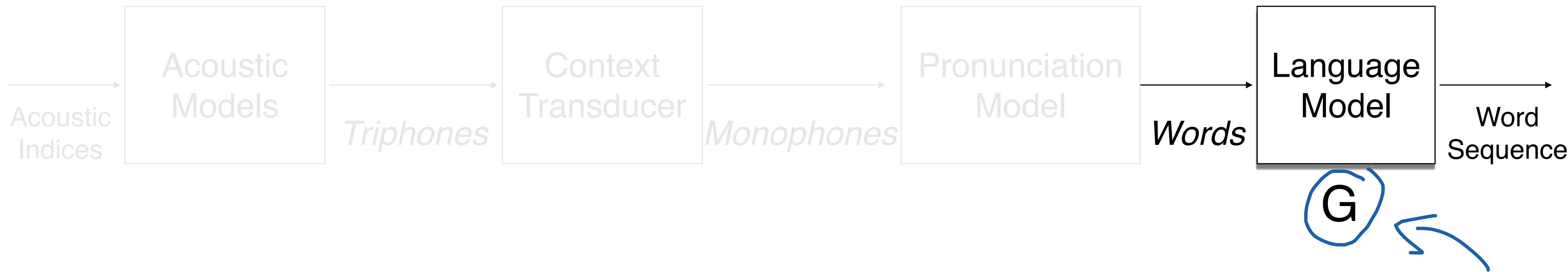
Determinized Dictionary WFST



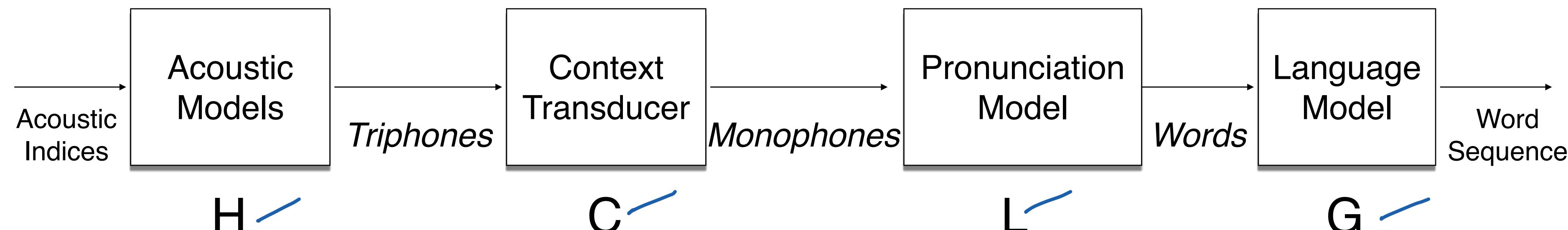
Minimized Dictionary WFST



WFST-based ASR System



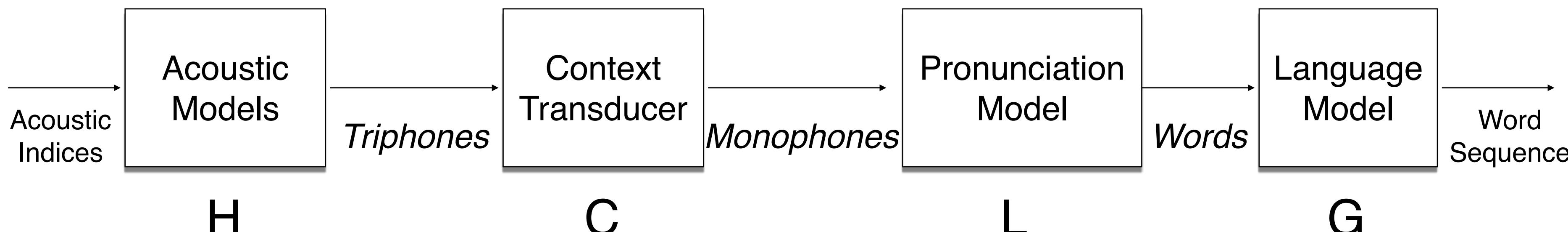
Decoding



Carefully construct a decoding graph D using optimization algorithms:

$$D = \min(\det(H) \circ \det(C) \circ \det(L) \circ G))$$

Decoding



Carefully construct a decoding graph D using optimization algorithms:

$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

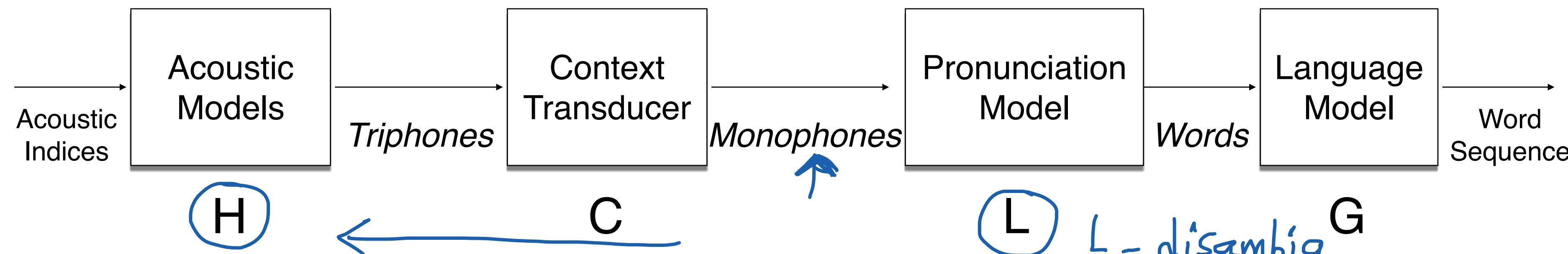
*homophones
E.g. read vs. reed
new vs. knew*

Need to add *disambiguation symbols* to L in order to make $L \circ G$ determinizable!
(auxiliary)

When are disambiguation symbols needed?

1. When you deal with homophones i.e., same pronunciation, different words.
2. When pronunciation of a word is a prefix of another. *cat caterpillar*

Decoding



Carefully construct a decoding graph D using optimization algorithms:

$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

\equiv

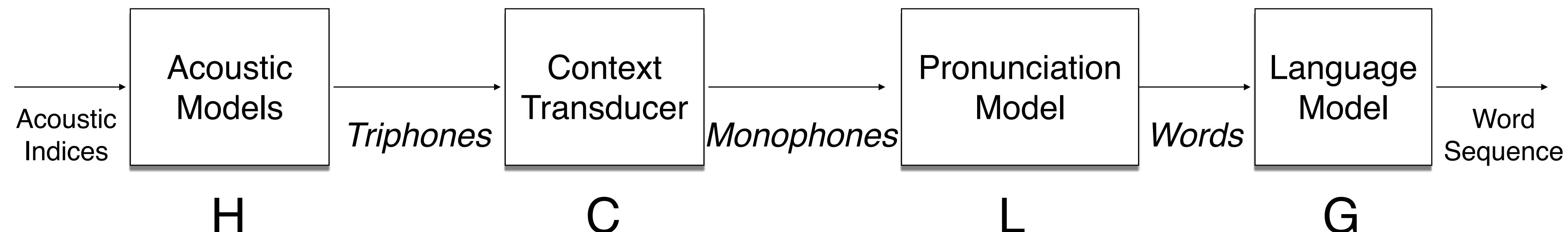
Need to add *disambiguation symbols* to L in order to make $L \circ G$ determinizable!

When are disambiguation symbols needed?

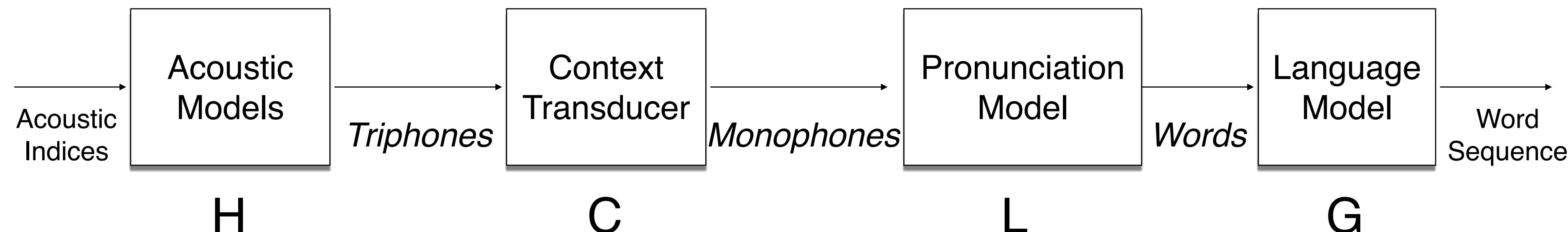
1. When you deal with homophones i.e., same pronunciation, different words.
2. When pronunciation of a word is a prefix of another.

book	b uh k #1
<u>books</u>	b uh k s
<u>right</u>	r ay t #1
<u>write</u>	r ay t #2
:	

Decoding



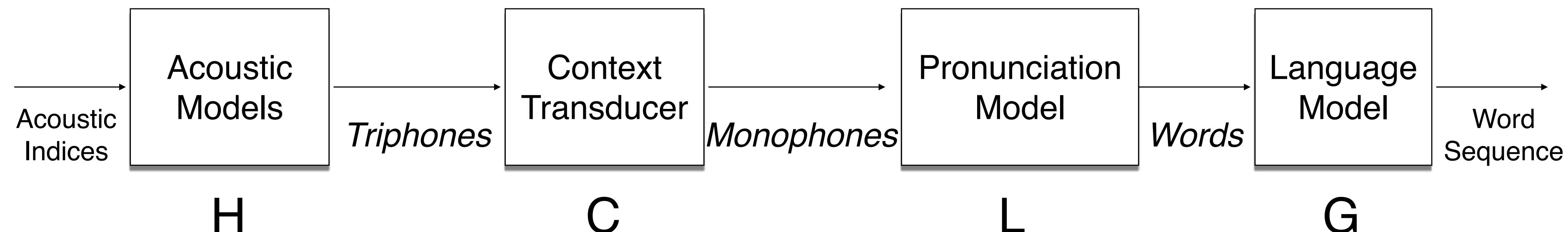
Decoding



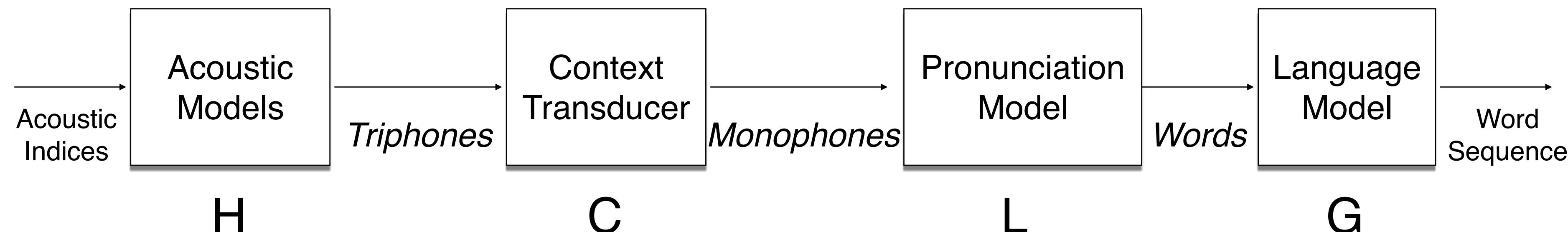
Carefully construct a decoding graph D using optimization algorithms:

$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Decoding



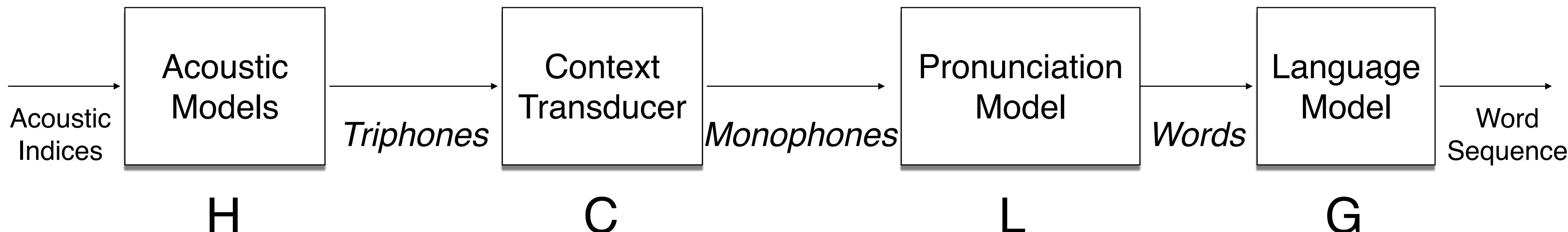
Decoding



Carefully construct a decoding graph D using optimization algorithms:

$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Decoding



Carefully construct a decoding graph D using optimization algorithms:

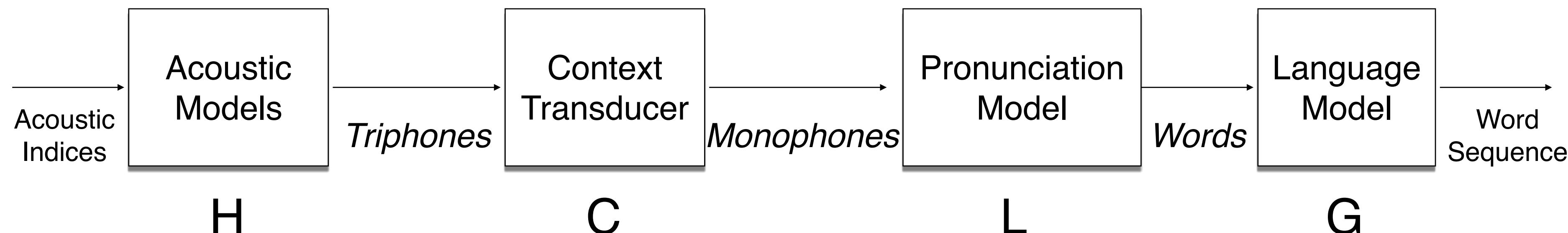
$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

O_1, \dots, O_T

Given a test utterance O , how do I decode it?

Assuming ample compute, first construct the following machine X from O .

Decoding

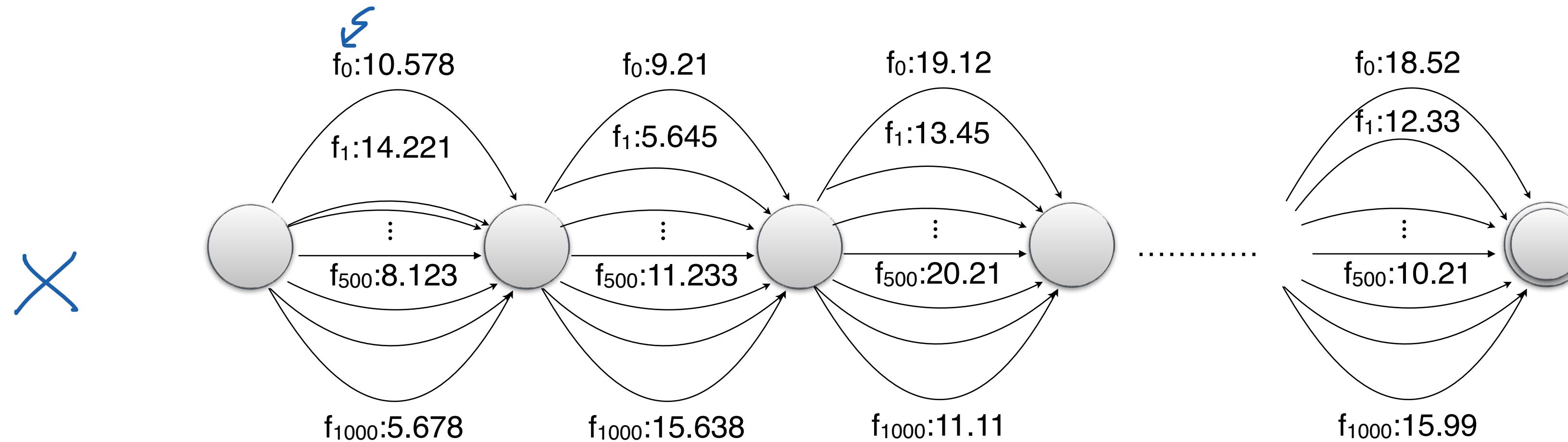


Carefully construct a decoding graph D using optimization algorithms:

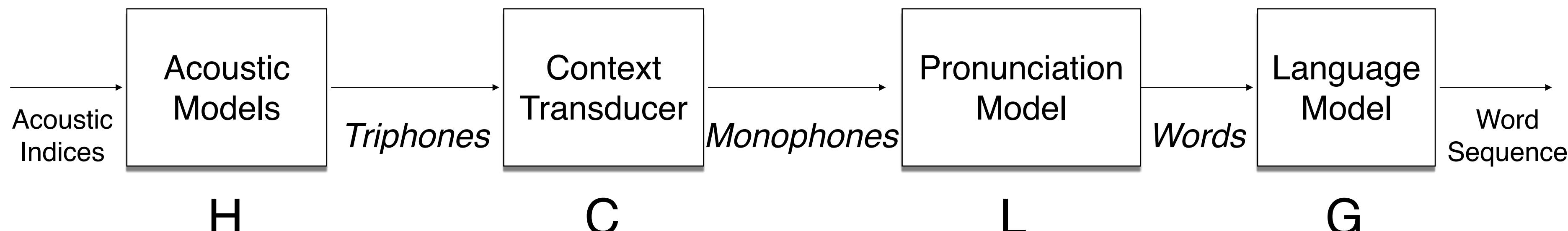
$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Given a test utterance O , how do I decode it?

Assuming ample compute, first construct the following machine X from O .



Decoding

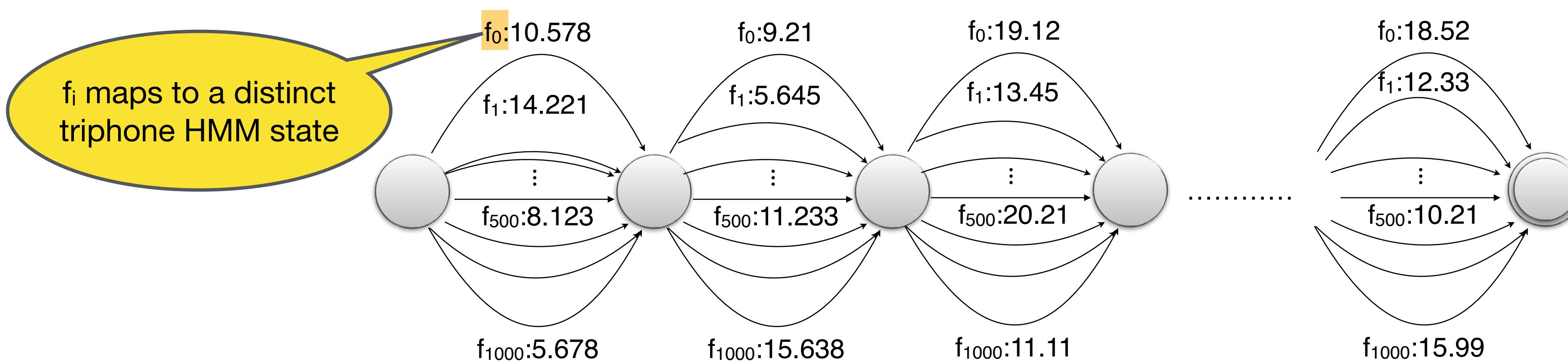


Carefully construct a decoding graph D using optimization algorithms:

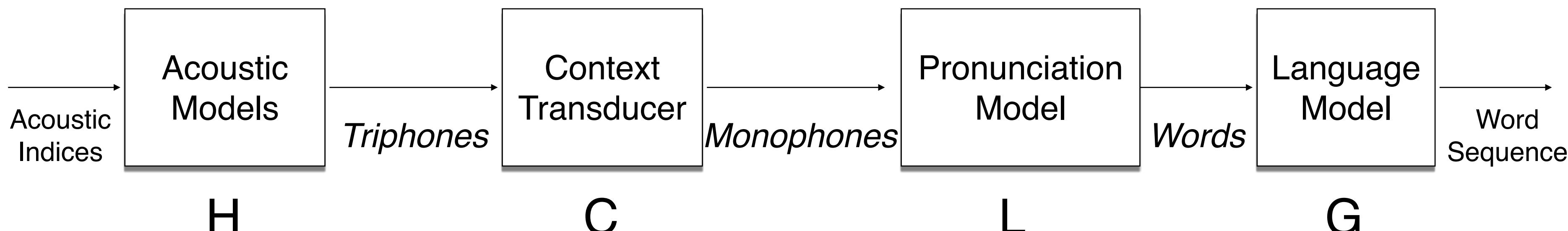
$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Given a test utterance O , how do I decode it?

Assuming ample compute, first construct the following machine X from O .



Decoding

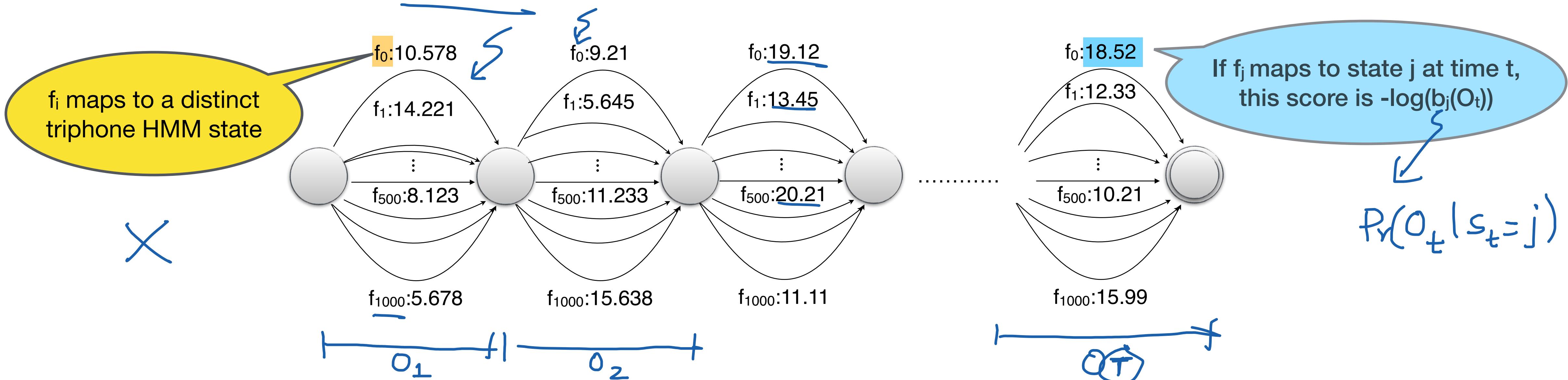


Carefully construct a decoding graph D using optimization algorithms:

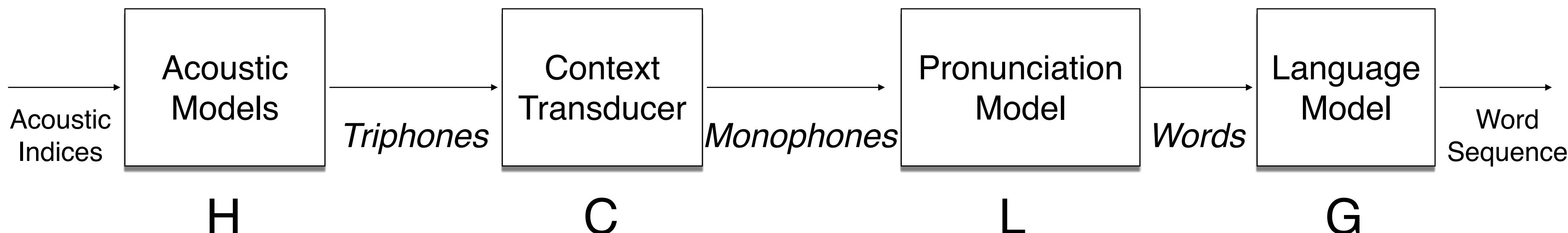
$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Given a test utterance O , how do I decode it?

Assuming ample compute, first construct the following machine X from O .



Decoding

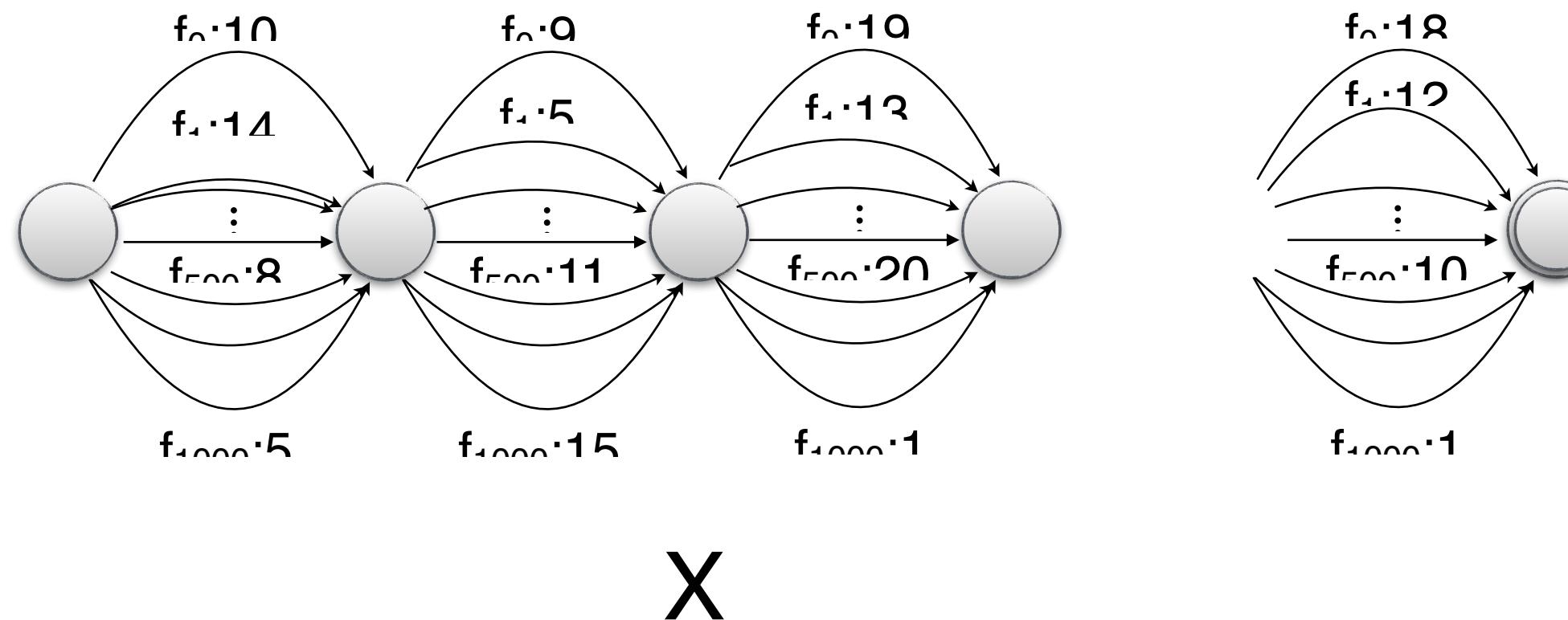


Carefully construct a decoding graph D using optimization algorithms:

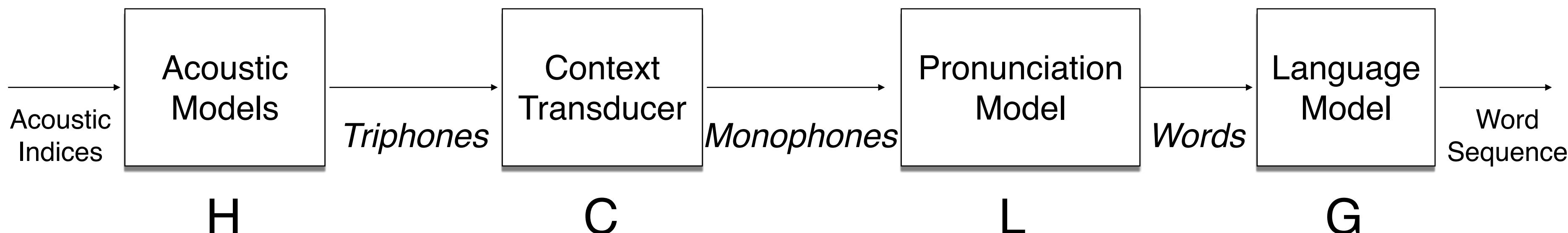
$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Given a test utterance O , how do I decode it?

Assuming ample compute, first construct the following machine X from O .



Decoding

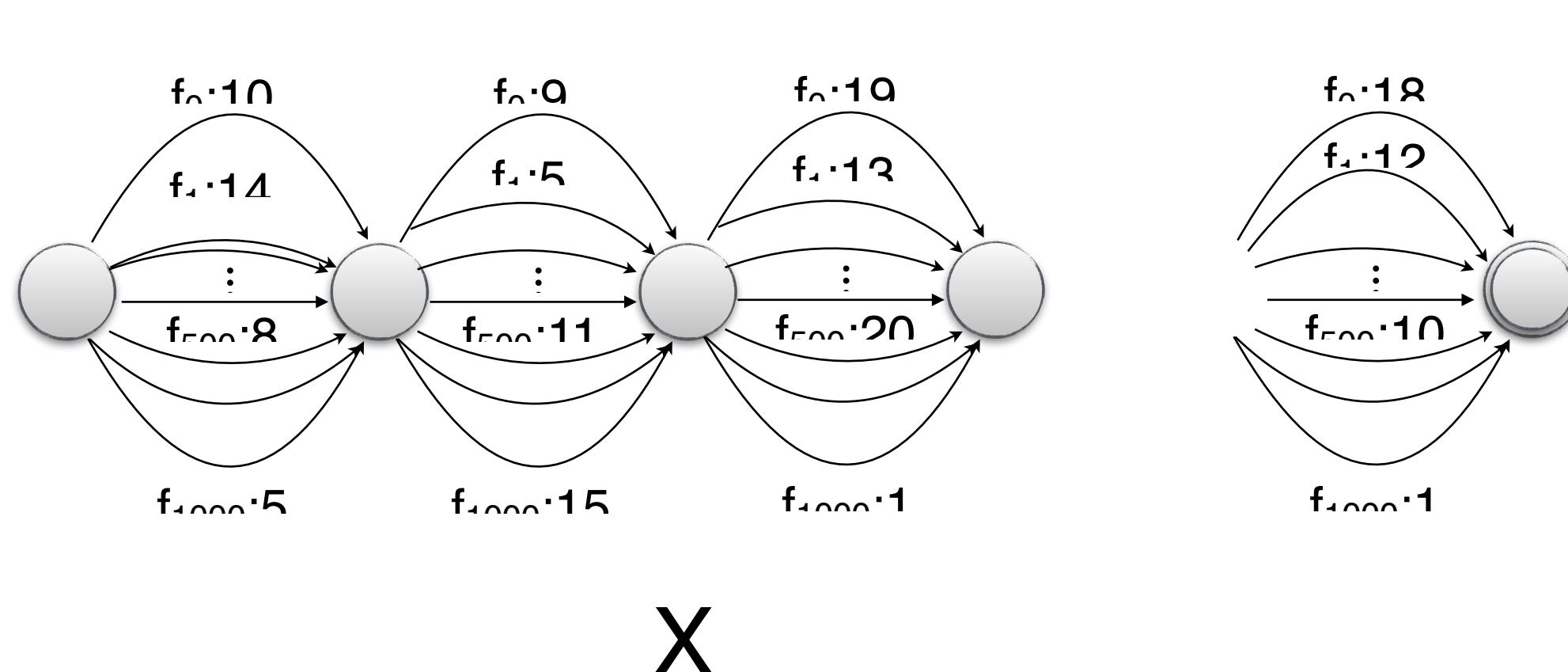


Carefully construct a decoding graph D using optimization algorithms:

$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Given a test utterance O , how do I decode it?

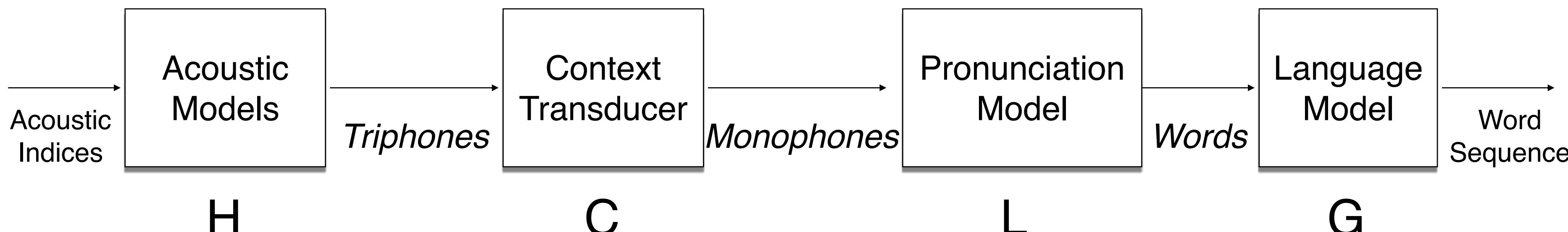
Assuming ample compute, first construct the following machine X from O .



$$W^* = \arg \min_{W=\text{out}[\pi]} \underbrace{X \circ D}_{\swarrow}$$

where π is a path in the composed FST
 $\text{out}[\pi]$ is the output label sequence of π

Decoding

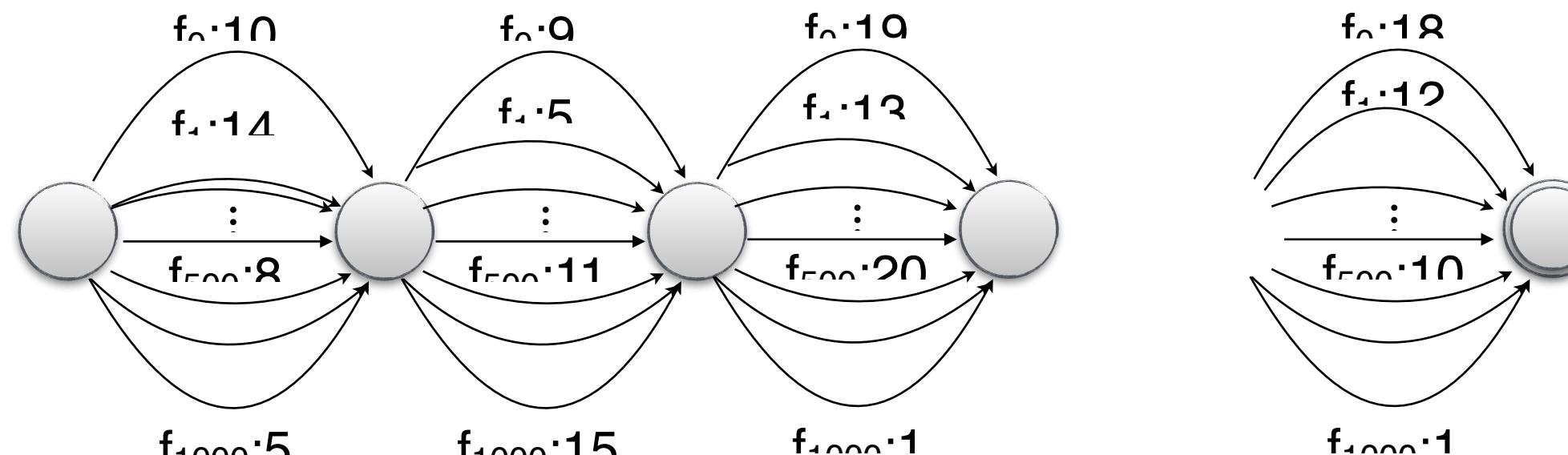


Carefully construct a decoding graph D using optimization algorithms:

$$D = \min(\det(H \circ \det(C \circ \det(L \circ G))))$$

Given a test utterance O , how do I decode it?

Assuming ample compute, first construct the following machine X from O .



X

$$W^* = \arg \min_{W=\text{out}[\pi]} X \circ D$$

where π is a path in the composed FST
 $\text{out}[\pi]$ is the output label sequence of π

X is never typically constructed;
 D is traversed dynamically using approximate search algorithms
 (discussed later in the semester)

Impact of WFST Optimizations

40K NAB Evaluation Set '95 (83% word accuracy)

network	states	transitions
G	1,339,664	3,926,010
$L \circ G$	8,606,729	11,406,721
$det(L \circ G)$	7,082,404	9,836,629
$C \circ det(L \circ G))$	7,273,035	10,201,269
$det(H \circ C \circ L \circ G)$	18,317,359	21,237,992

Impact of WFST Optimizations

40K NAB Evaluation Set '95 (83% word accuracy)

network	states	transitions
G	1,339,664	3,926,010
$L \circ G$	8,606,729	11,406,721
$det(L \circ G)$	7,082,404	9,836,629
$C \circ det(L \circ G))$	7,273,035	10,201,269
$det(H \circ C \circ L \circ G)$	18,317,359	21,237,992

network	x real-time
$C \circ L \circ G$	<u>12.5</u>
$C \circ det(L \circ G)$	<u>1.2</u>
$det(H \circ C \circ L \circ G)$	<u>1.0</u>
$push(min(F))$	0.7