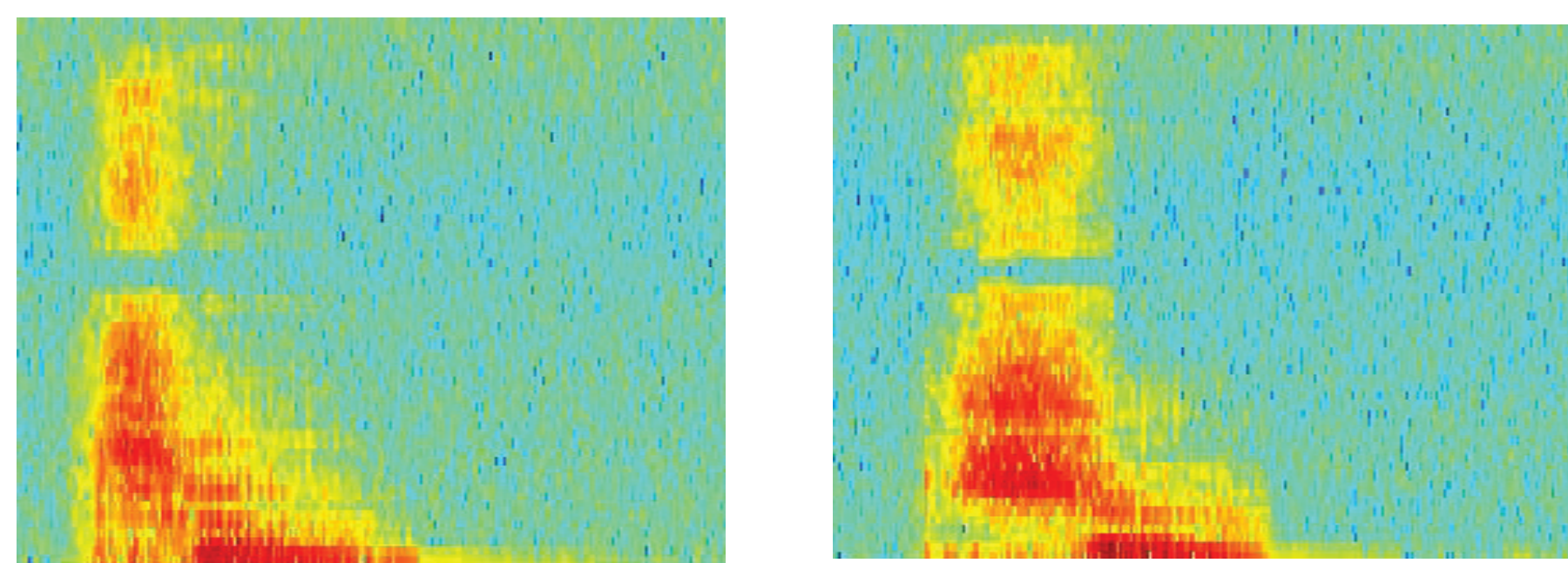


SPREAD

SPREAD is a novel agent-based sound perception model. It simulates how sound features are propagated, attenuated, spread and degraded as they traverse a virtual environment.

Problem

SPREAD has been demonstrated to work with 100 different environment sounds and can be the basis of agent-environment interactions. However environment sounds are not enough to simulate agent-to-agent or user-to-agent interactions.



Less differences between phonemes provide a bigger challenge in the packet generation and matching stage (left: "t"; right: "ch").

Goal: A user should be able call the name of an agent through speech, and the agent will either respond or not respond based on their location in the scene.

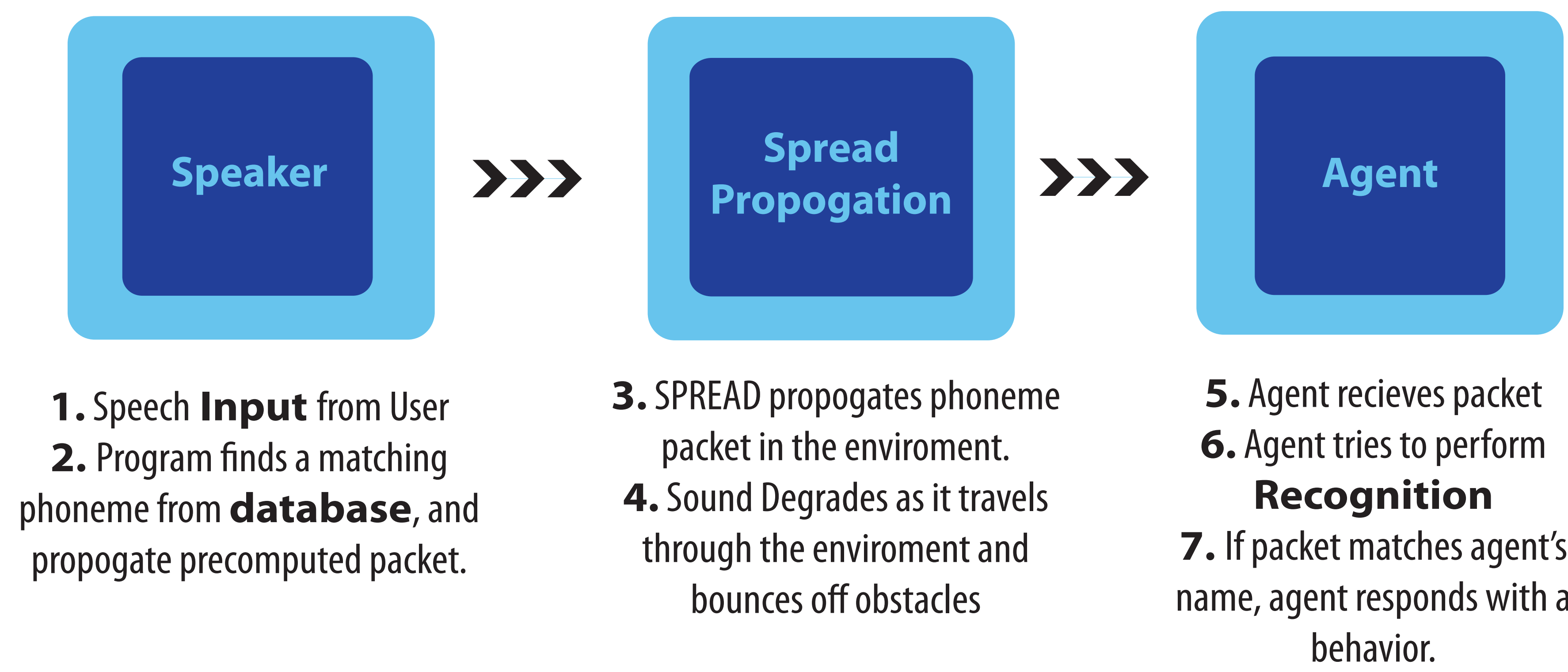
Phonemes-Approach

Spoken speech can be broken down. The smallest unit linguistically is the "phoneme. The English language consists of about 44 different phonemes.


20 Vowel Phonemes for American English			
AKSES Character	Phoneme Name	AKSES Character	Phoneme Name
a_A	aa	a_A	AA
a_O	ao	a_O	AO
e_E	ee	e_E	EE
e_I	ei	e_I	EI
i_I	ii	i_I	II
o_O	oo	o_O	OO
o_U	ou	o_U	OU
u_U	uu	u_U	UU
sch	sch	sch	SCH
ph	ph	ph	PH
ng	ng	ng	NG
kw	kw	kw	KW
gh	gh	gh	GH
sh	sh	sh	SH
zh	zh	zh	ZH
ch	ch	ch	CH
ts	ts	ts	TS
th	th	th	TH
ns	ns	ns	NS
lh	lh	lh	LH
ms	ms	ms	MS
ks	ks	ks	KS
ps	ps	ps	PS
kw	kw	kw	KW
gh	gh	gh	GH
sh	sh	sh	SH
zh	zh	zh	ZH
ch	ch	ch	CH
ts	ts	ts	TS
th	th	th	TH
ns	ns	ns	NS
lh	lh	lh	LH
ms	ms	ms	MS
ks	ks	ks	KS
ps	ps	ps	PS

Each agent is given the shortest name possible -- one phoneme long. This makes the pipeline easy to analyze and debug.

Pipeline



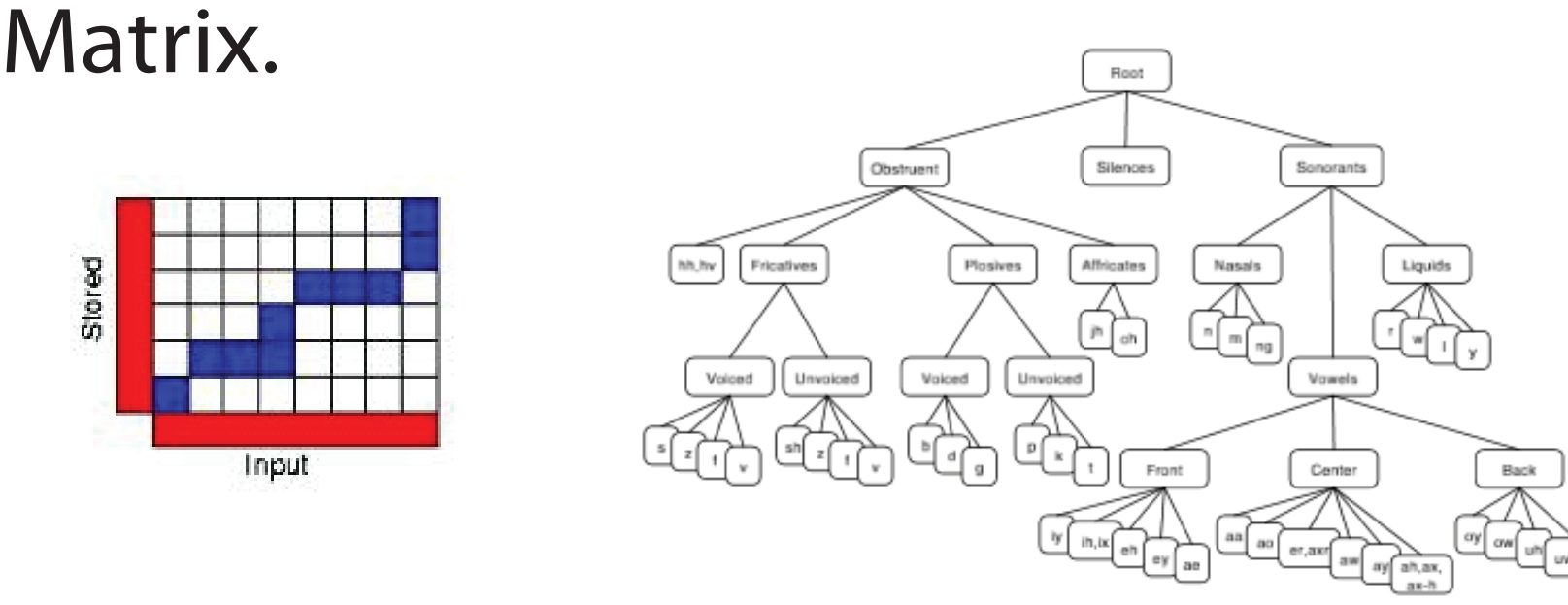
Input



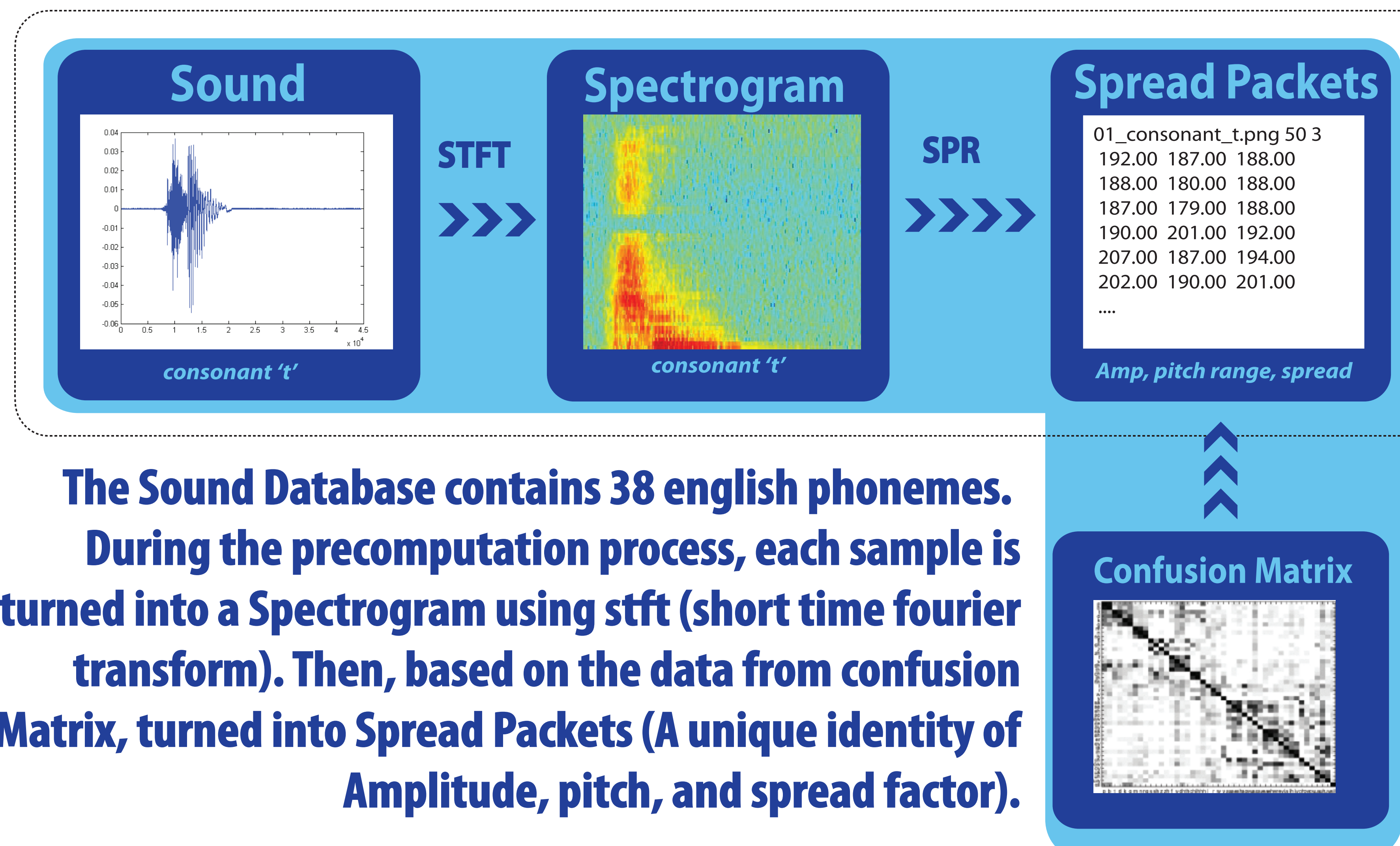
Speech Input from user is first processed through **Speech SDK 5.1 for phonemes.**

Recognition

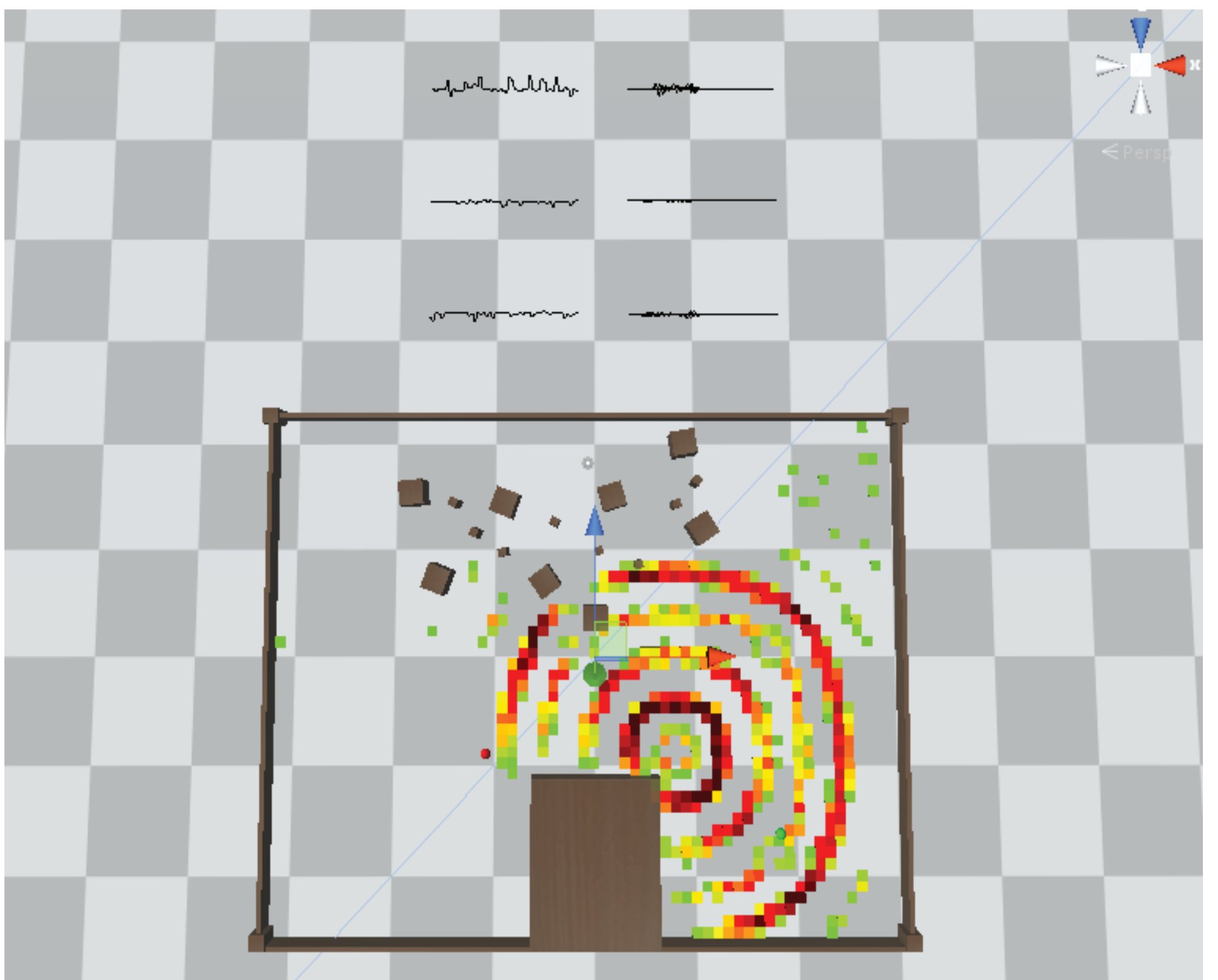
Degraded Packets are first matched to original packets via **Dynamic Time Warping**. Confidence scores for matches are calculated using an **HCA tree** calculated from the Similarity Matrix.



Sound Database



Evaluation



Three agents with the names "t", "ah", and "m" are located in the scene. each agent was able to recognize their name when they are called. (In the screen shot, the propagated sound is "t". The agent responds by turning red).

However for some phonemes, analysis of second or third phoneme matches makes very little sense. (ie: when propagating "t", the second match will be the "ei" *eight)

Future Work

- Sample using a logarithmic scheme (instead of a uniform scheme) to better simulate the biological perception model.
- Implement multiple-phoneme propagation and recognition.
- Introduce into the system, a statistics based guessing scheme that will predict the next phoneme in a sequence.