

# DeepSpeech: A Scalable Decoding System that Integrates Knowledge for Speech Recognition

Zifei Shan, Tianxin Zhao, Haowen Cao  
{zifei, tianxin, caohw}@stanford.edu

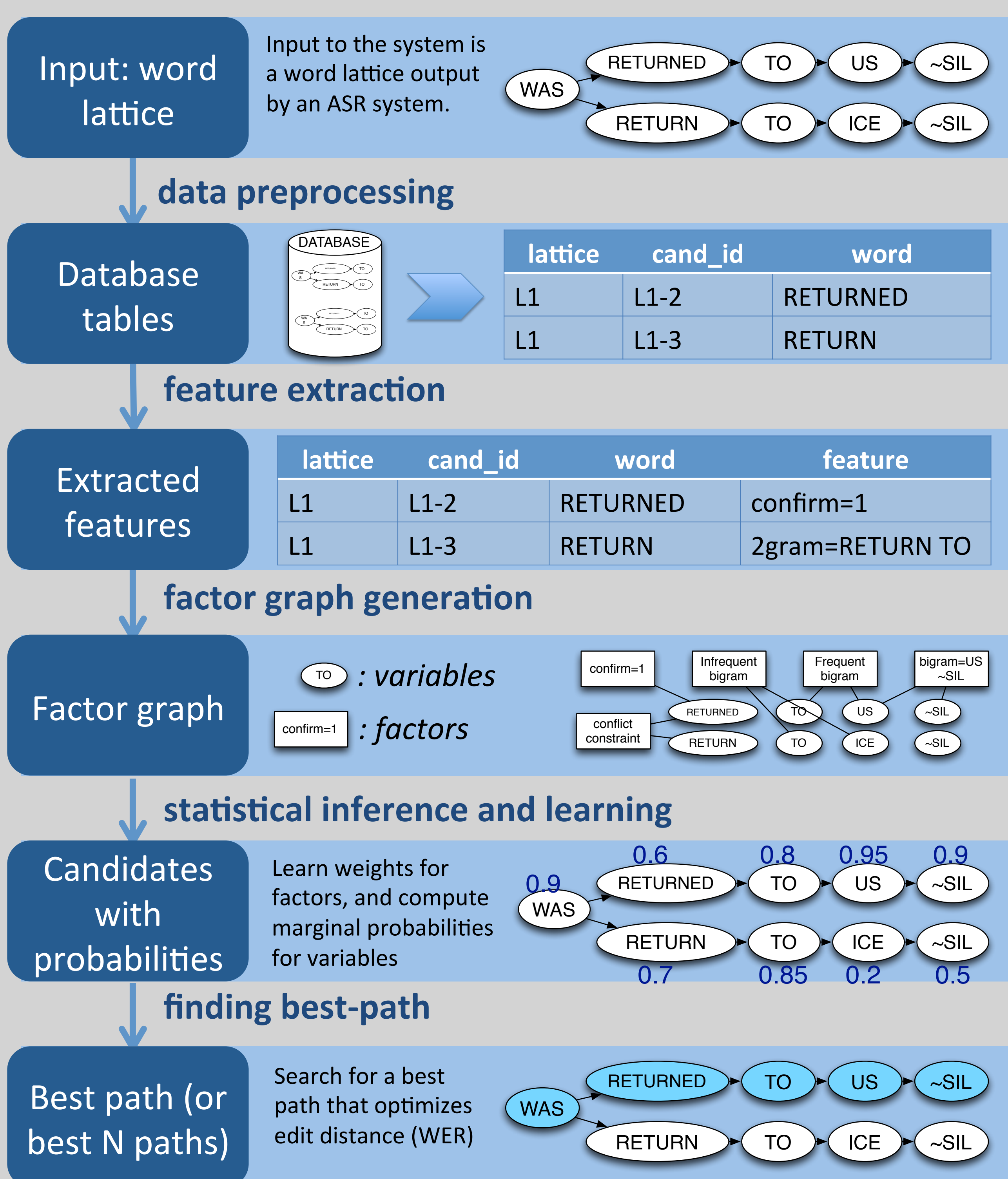
## Executive Summary

- DeepSpeech flexibly integrates **different levels of knowledge** to decode a word lattice in speech recognition within a word-level CRF model, in an interpretable manner.
- DeepSpeech facilitates **feature extraction, factor graph generation, and statistical learning and inference**. It takes word lattice as input, perform feature extraction specified by developers, generate factor graphs based on descriptive rules, and perform learning and inference automatically.
- DeepSpeech is based on the scalable statistical inference engine **DeepDive** (<http://deepdive.stanford.edu>).

## DeepSpeech Overview

<b>Problem</b>	How to jointly integrate different levels of knowledge in speech recognition?
<b>Solution</b>	Decoding based on Conditional Random Fields that integrates various features.
<b>Results</b>	Got WER <b>10.2%</b> on 150k broadcast news lattices in <b>near real-time</b> , with a simple feature set. ( <i>baseline 22.9%, oracle 2.1%</i> )
<b>Future</b>	Candidate generation with linguistic knowledge might beat oracle error rate; joint inference on acoustic and language models

## The Architecture of DeepSpeech

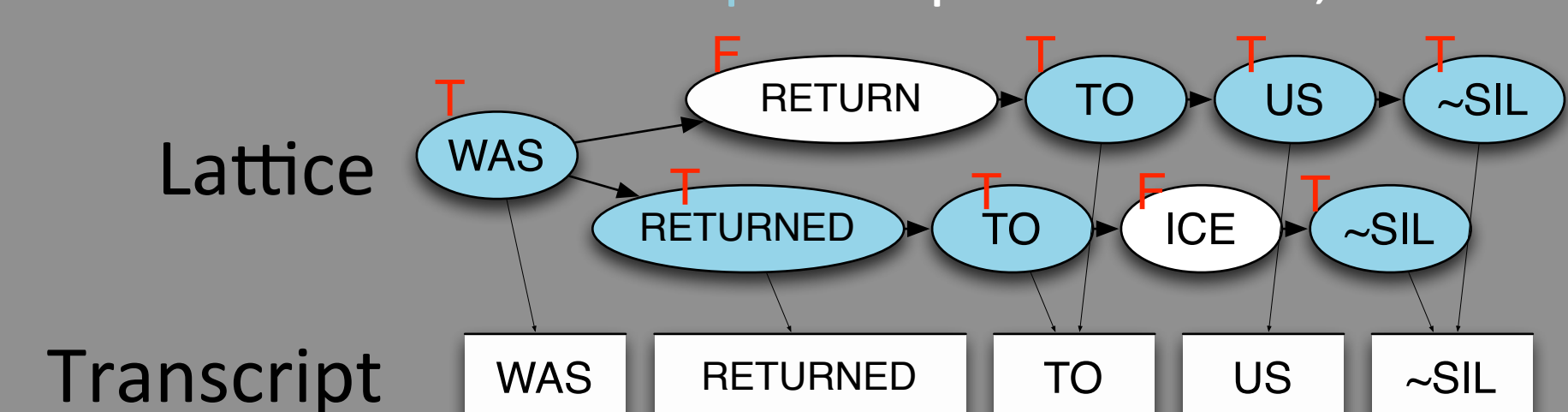


## Experimental Setup & Preliminary Results

Features	Datasets																												
After initial feature engineering, our current system implements following features:	We train and test on broadcast news lattices (LDC2011T06, 150k lattices). We holdout 50% of training set for testing.																												
1. Unigram and bigram frequency in Google Ngram. (skip "silence")	<b>Results</b> We use SCLITE for scoring. We evaluate a baseline system (Attlia), DeepSpeech and lattice oracle (optimal) error rate.																												
2. All bigrams around "silence"																													
3. POS tag 2gram and 3gram																													
4. Candidates that overlap in time cannot be both true ( <i>a CRF rule that indicates constraint</i> )																													
5. Candidates on a same path should be true at same time ( <i>a linear-chain CRF rule</i> )																													
Next steps: co-reference features and candidate generation	<table border="1"> <thead> <tr> <th>System</th> <th>Corr</th> <th>Sub</th> <th>Del</th> <th>Ins</th> <th>Err</th> <th>S.Err</th> </tr> </thead> <tbody> <tr> <td>Baseline</td> <td>77.8</td> <td>5.4</td> <td>16.8</td> <td>0.6</td> <td>22.9</td> <td>96.9</td> </tr> <tr> <td>DeepSpeech</td> <td>92.0</td> <td>3.6</td> <td>4.3</td> <td>2.2</td> <td><b>10.2</b></td> <td>75.7</td> </tr> <tr> <td>Oracle</td> <td>99.9</td> <td>0.0</td> <td>0.1</td> <td>2.0</td> <td>2.1</td> <td>50.8</td> </tr> </tbody> </table>	System	Corr	Sub	Del	Ins	Err	S.Err	Baseline	77.8	5.4	16.8	0.6	22.9	96.9	DeepSpeech	92.0	3.6	4.3	2.2	<b>10.2</b>	75.7	Oracle	99.9	0.0	0.1	2.0	2.1	50.8
System	Corr	Sub	Del	Ins	Err	S.Err																							
Baseline	77.8	5.4	16.8	0.6	22.9	96.9																							
DeepSpeech	92.0	3.6	4.3	2.2	<b>10.2</b>	75.7																							
Oracle	99.9	0.0	0.1	2.0	2.1	50.8																							
	<b>Performance:</b> DeepSpeech runs <b>70min</b> for training and testing, while this dataset is <b>~400 hours</b> of speech (real-time!)																												

### Distant supervision to obtain training data

- We use distant supervision techniques to get training labels on candidate level: given a lattice and its transcript, we:
- Find optimal paths in the lattice that matches the transcript with Dynamic Programming
  - Label all **matched** words in all **optimal** paths as **true**, others as false



## DeepSpeech has three features!

### Extract & Integrate linguistic features

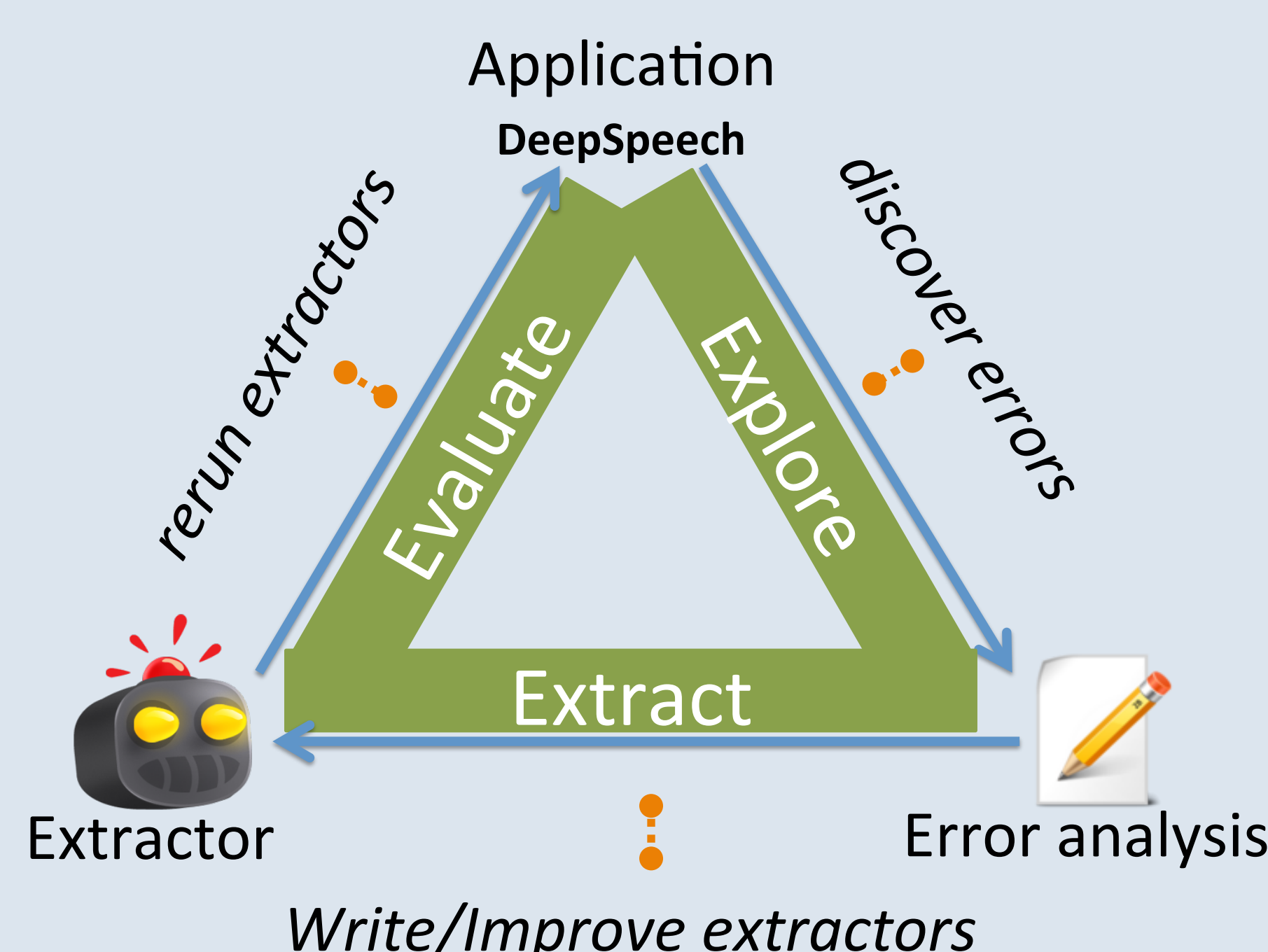
DeepSpeech provides a framework for developers to extract and integrate high-level features

Here is how developers can easily plug-in a coreference feature "extractor" to DeepSpeech:

SQL	Python
Define a SQL query to generate all pairs of candidate words that appear in the same lattice, and pair it with a python function.	We write a Python function to process all phrase pairs and identify coreferent pairs.
<pre>SELECT t0.CID, t0.TEXT,        t1.CID, t1.TEXT FROM candidate t0,      candidate t1 WHERE t0.LID = t1.LID USEPYTHON pyfunc</pre>	<pre>def pyfunc(c1, t1, c2, t2):     if edit_dist(t1, t2) &lt; 2:         emit("Coref", c1, c2)</pre>

### Simpler Feature Engineering

DeepSpeech supports an "E3 loop" for feature engineering



### Rigorous Probabilistic Framework

DeepSpeech uses a joint probability model that enables rigorous probabilistic interpretation

