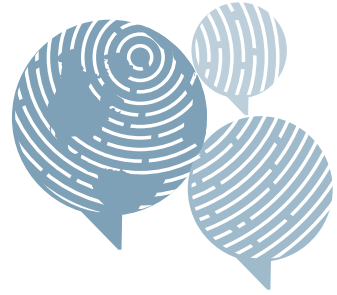# Improvements to Urdu-English
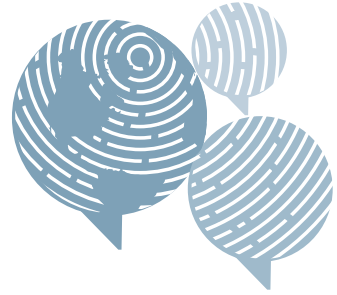
# SCALE Summer Workshop Results

## Chris Callison-Burch

Kathy Baker, Steven Bethard, Michael Bloodgood, Ralf Brown, Glen Coppersmith, Bonnie Dorr, Wes Filardo, Kendall Giles, Ann Irvine, Mike Kayser, Lori Levin, Justin Martineau, Jim Mayfield, Scott Miller, Aaron Phillips, Andrew Philpot, Christine Piatko, Lane Schwartz, and David Zajic

# human language technology
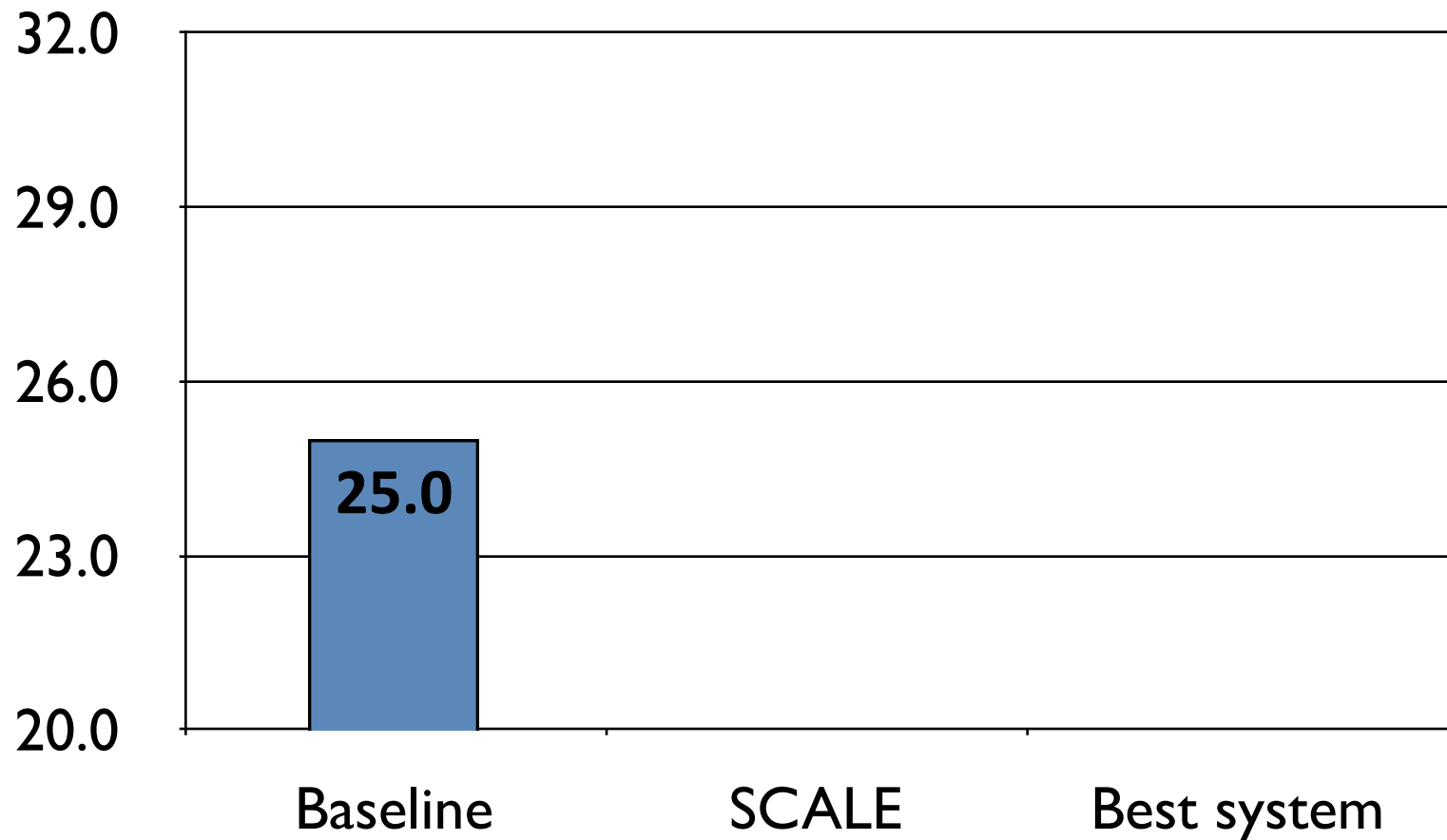c e n t e r   o f   e x c e l l e n c e
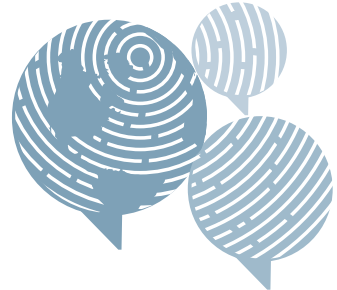
# The Punch Line

# The Punch Line
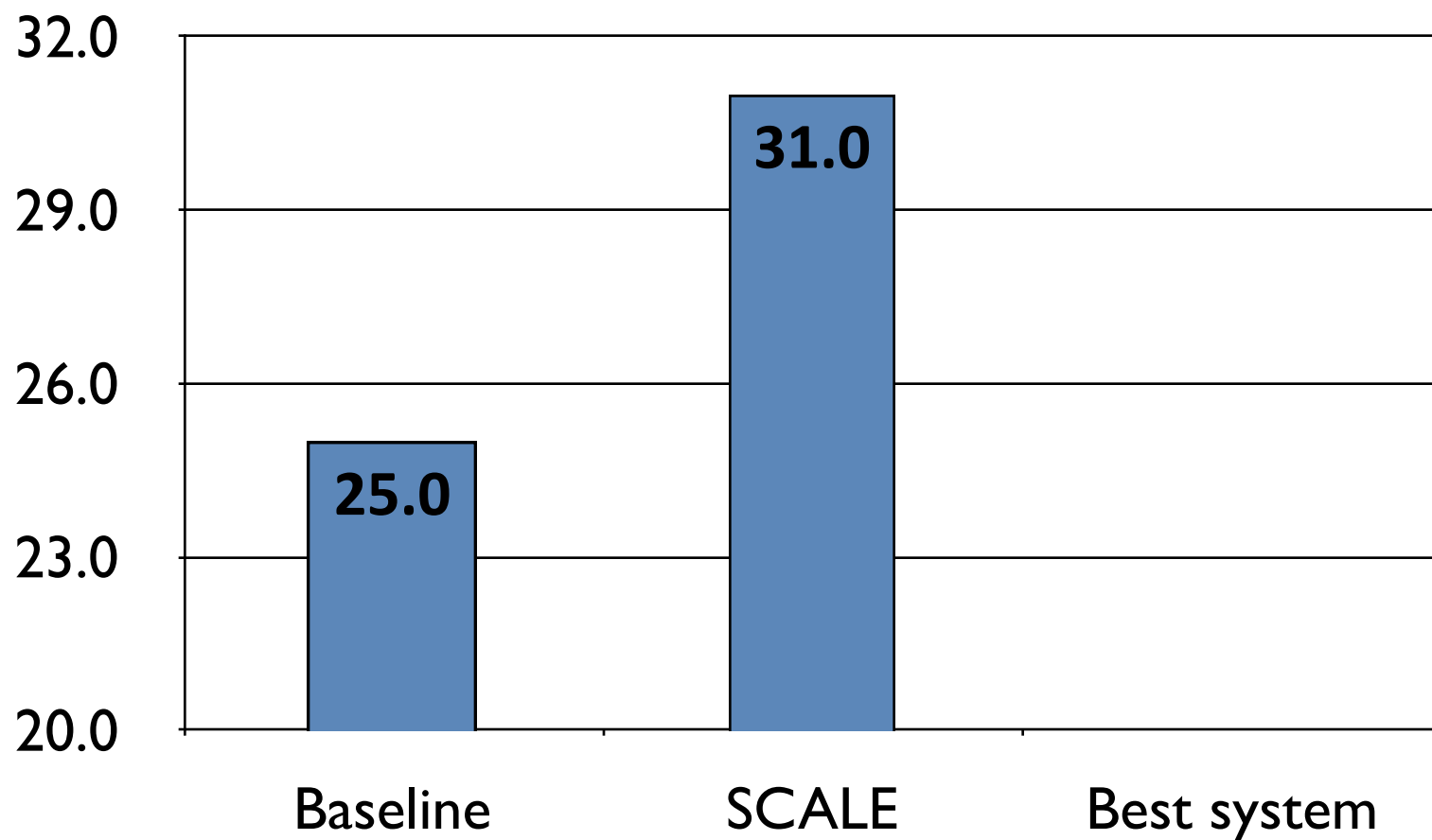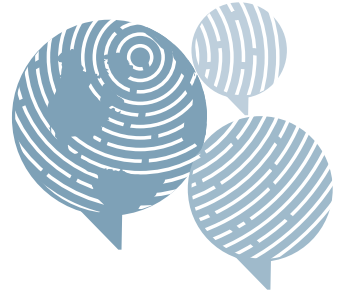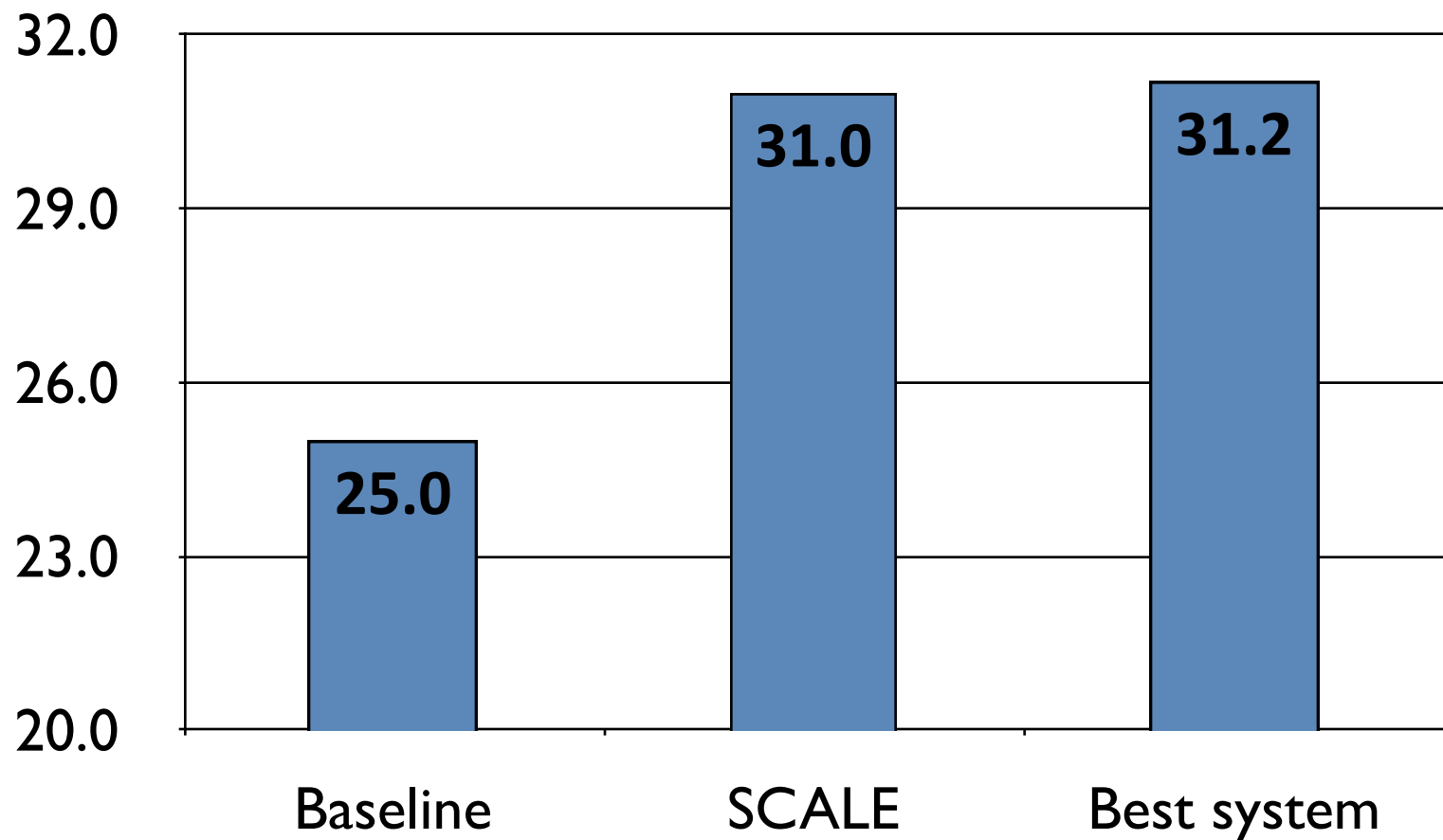
Bleu score on blind NIST Urdu-English test set

# The Punch Line

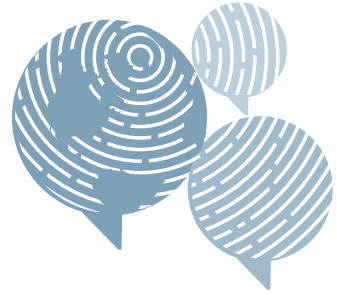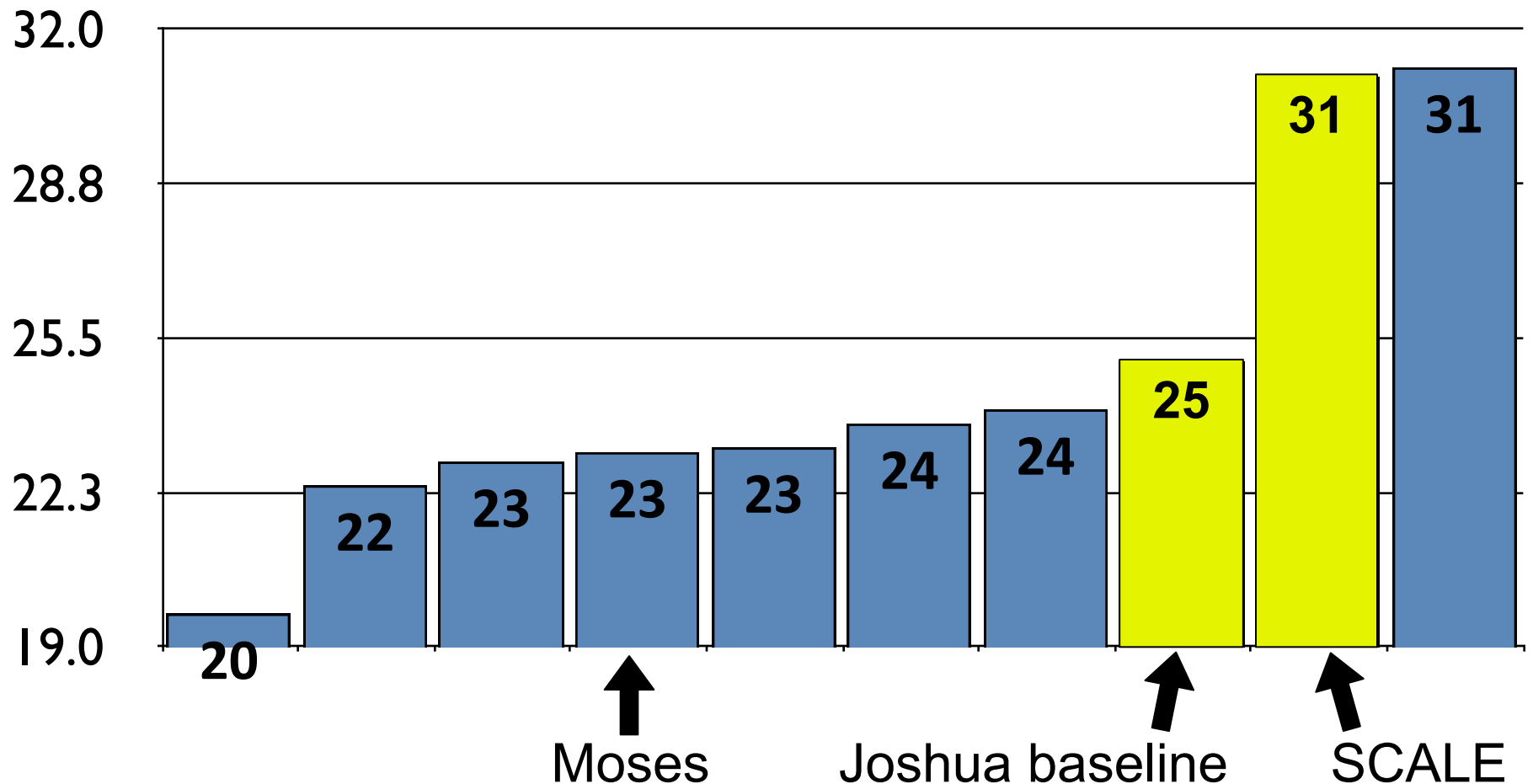Bleu score on blind NIST Urdu-English test set

# The Punch Line

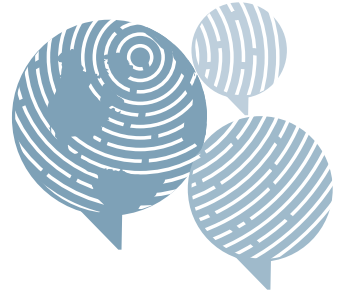Bleu score on blind NIST Urdu-English test set
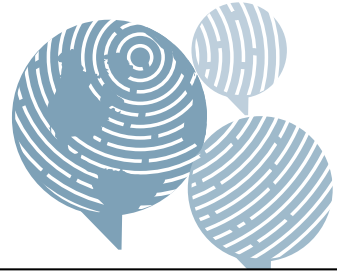
# The field for the Urdu task

All system scores on NIST09 Urdu-English constrained task

# How we got the gains

- Improvements over baseline decoder
  - Syntax-based model replaced hierarchal phrase-based model
  - Semantic elements (named entities, modalities) are grafted onto syntax trees
  - Unknown Urdu words are transliterated instead of left untranslated
  - Larger number of feature fns score translations
  - Word alignments improved by incorporating multiple aligners including a syntax-based aligner

# Translation improvements

**'first nuclear experiment in 1990 was'**
Thomas red Unilever National Laboratory of the United States وپن in designer, are already working on the book of Los ايلموس National Laboratory ڈینی, former director of the technical انٹیلجنس written with the cooperation of سٹلمین.
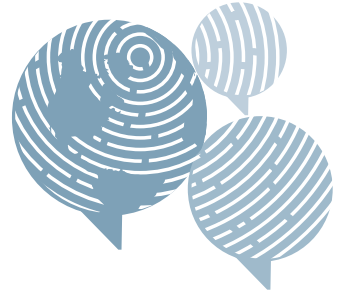
This book 'nuclear express: political history and the expansion of bomb' has been written, and the two writers have also claimed that the country has made nuclear bomb is he or any other country's nuclear secrets to or that of any other nuclear چرائے power cooperation is achieved.

**The First Nuclear Test Was in 1990.**

Thomas red of the United States, the National Laboratory in designer are already working on the book of Los Alamos National Laboratory, former director of the technical intelligence, with the cooperation of Diana steelman wrote.

This book under the title of the spread of nuclear expressway: the political history of the bomb and this has been written and the two writers have claimed that the country also has made nuclear bomb or any other country, Korea nuclear secrets, or any of the other nuclear power cooperation.

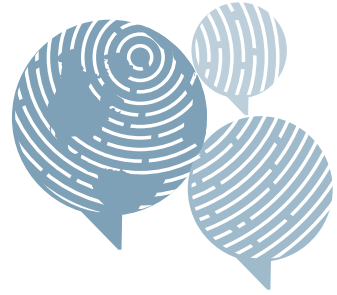# Who did what to whom?

**Baseline**

He said that China, North Korea, Iran, Syria, Pakistan, through Egypt, Libya and Yemen is to provide nuclear technology.

Thomas was red when this question why China has provided the nuclear technology to Pakistan, In response, He said as China and India was joint enemy of Pakistan.
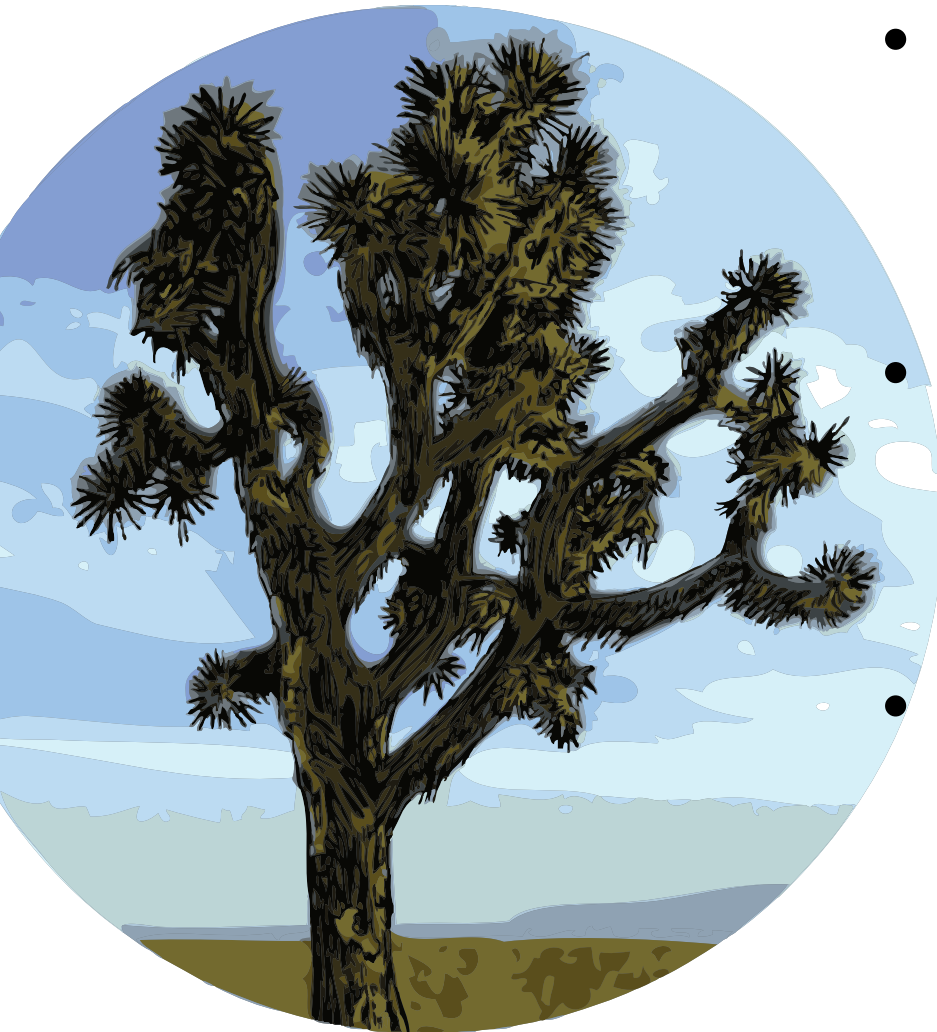
**SCALE final system**

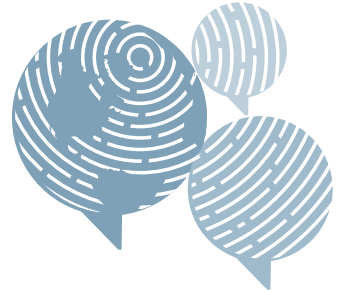He said that China would provide nuclear technology to North Korea, Iran, Syria, Pakistan, Egypt, Libya and Yemen.

Thomas red when was this question why China has provided to Pakistan nuclear technology, he said in response to China, Pakistan and India as a common enemy.

# Joshua Decoder

- Synchronous context free grammars generate pairs of corresponding strings

- Can be used to describe translation and re-ordering between languages

- Because Joshua uses SCFGs, it translates sentences by parsing them

# Example SCFG for Urdu

|  | Urdu | English |
|---|---|---|
| S → | NP① VP② | NP① VP② |
| VP→ | PP① VP② | VP② PP① |
| VP→ | V① AUX② | AUX② V① |
| PP → | NP① P② | P② NP① |
| NP → | *hamd ansary* | *Hamid Ansari* |
| NP → | *na}b sdr* | *Vice President* |
| V → | *namzd* | *nominated* |
| P → | *kylye* | *for* |
| AUX → | *taa* | *was* |

*hamd ansary    na}b sdr    kylye    namzd    taa*

NP❶

hamd ansary     na}b sdr     kylye     namzd     taa

NP❶

Hamid Ansari

NP❶

hamd ansary

NP❷

na}b sdr     kylye     namzd     taa

NP❶

Hamid Ansari

NP❷

Vice President

NP❶

hamd ansary

NP❷

na}b sdr

P❸

kylye     namzd     taa

NP❶

Hamid Ansari
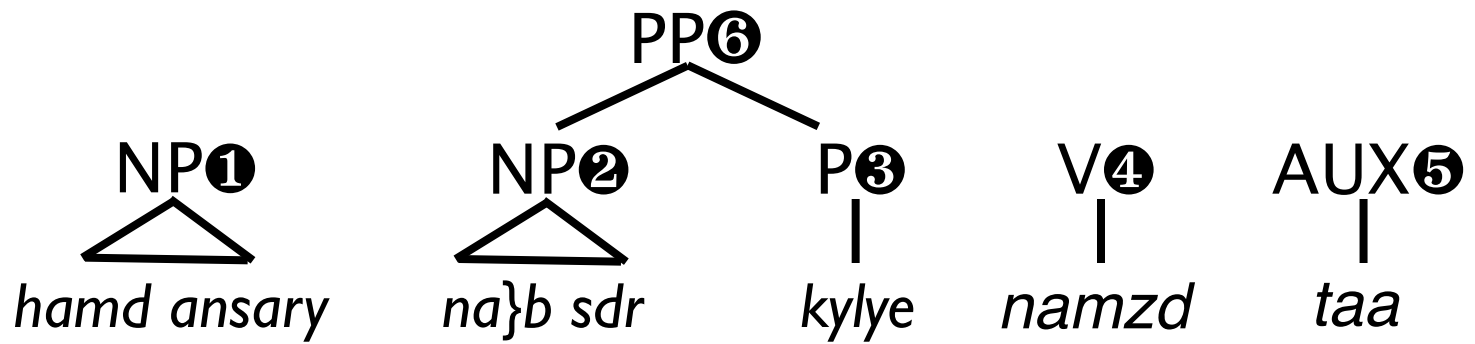
NP❷

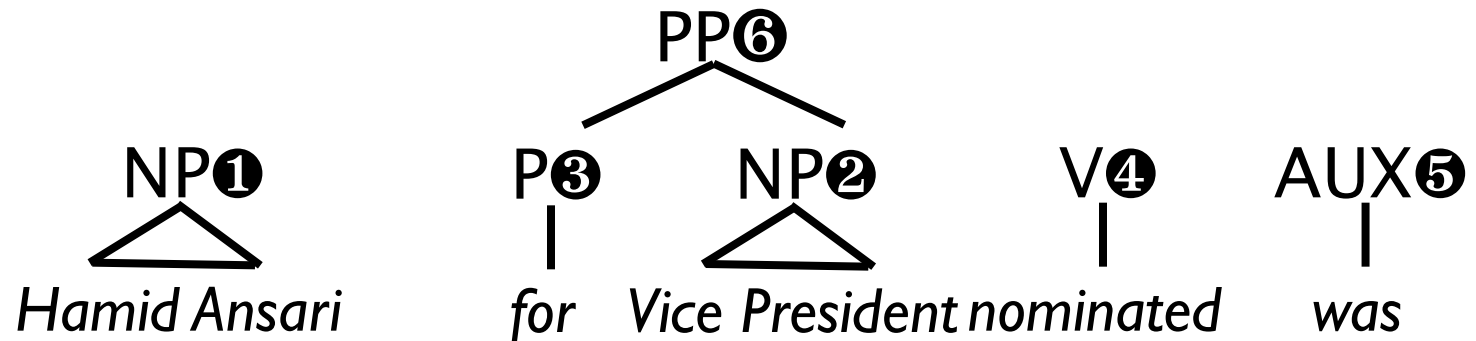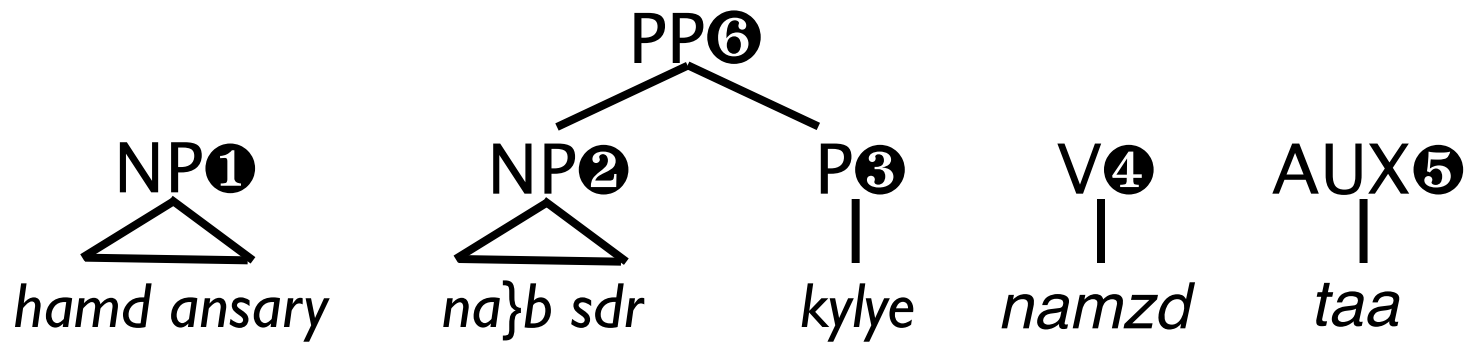Vice President

P❸

for

NP❶
*hamd ansary*

NP❷
*na}b sdr*

P❸
*kylye*

V❹
*namzd*

*taa*

NP❶
*Hamid Ansari*

NP❷
*Vice President*

P❸
*for*

V❹
*nominated*

NP❶     NP❷     P❸     V❹     AUX❺

*hamd ansary*    *na}b sdr*    *kylye*    *namzd*    *taa*

NP❶     NP❷     P❸     V❹     AUX❺

*Hamid Ansari*    *Vice President*    *for*    *nominated*    *was*

NP❶ NP❷ P❸ V❹ AUX❺

*hamd ansary*  *na}b sdr*  *kylye*  *namzd*  *taa*

NP❶ NP❷ P❸ V❹ AUX❺

*Hamid Ansari*  *Vice President*  *for*  *nominated*  *was*

PP❻

NP❶   NP❷   P❸   V❹   AUX❺

*hamd ansary*   *na}b sdr*   *kylye*   *namzd*   *taa*

PP❻

NP❶   P❸   NP❷   V❹   AUX❺

*Hamid Ansari*   *for*   *Vice President*   *nominated*   *was*

PP❻

NP❶  NP❷  P❸  V❹  AUX❺

*hamd ansary*  *na}b sdr*  *kylye*  *namzd*  *taa*

PP❻

NP❶  P❸  NP❷  V❹  AUX❺

*Hamid Ansari*  *for*  *Vice President*  *nominated*  *was*

**PP❻**

**NP❶**

NP❷    P❸

*hamd ansary*    *na}b sdr*    *kylye*

**V❹**    **AUX❺**

*namzd*    *taa*

**VP❼**

---

**PP❻**

**NP❶**

P❸    NP❷

*Hamid Ansari*    *for*    *Vice President*

**AUX❺**    **V❹**

*was*    *nominated*

**VP❼**

**Top tree:**

NP❶: *hamd ansary*

PP❻
- NP❷: *na}b sdr*
- P❸: *kylye*

VP❼
- V❹: *namzd*
- AUX❺: *taa*

**Bottom tree:**

NP❶: *Hamid Ansari*

PP❻
- P❸: *for*
- NP❷: *Vice President*

VP❼
- AUX❺: *was*
- V❹: *nominated*

NP❶ *hamd ansary*  NP❷ *na}b sdr*  P❸ *kylye*  V❹ *namzd*  AUX❺ *taa*

PP❻  VP❼  VP❽

NP❶ *Hamid Ansari*  AUX❺ *was*  V❹ *nominated*  P❸ *for*  NP❷ *Vice President*

VP❼  PP❻  VP❽

## Tree 1

- VP❽
  - PP❻
    - NP❶ — *hamd ansary*
    - NP❷ — *na}b sdr*
    - P❸ — *kylye*
  - VP❼
    - V❹ — *namzd*
    - AUX❺ — *taa*

## Tree 2

- VP❽
  - VP❼
    - NP❶ — *Hamid Ansari*
    - AUX❺ — *was*
    - V❹ — *nominated*
  - PP❻
    - P❸ — *for*
    - NP❷ — *Vice President*

## Tree 1

- S❾
  - NP❶ — *hamd ansary*
  - VP❽
    - PP❻
      - NP❷ — *na}b sdr*
      - P❸ — *kylye*
    - VP❼
      - V❹ — *namzd*
      - AUX❺ — *taa*

## Tree 2

- S❾
  - NP❶ — *Hamid Ansari*
  - VP❽
    - VP❼
      - AUX❺ — *was*
      - V❹ — *nominated*
    - PP❻
      - P❸ — *for*
      - NP❷ — *Vice President*

# Hiero-style SCFG rules

- Baseline only supported Hiero-style SCFG rules which have only one non-terminal symbol

- Not as nice as linguistically motivated rules, do not capture the reordering in Urdu

$X_1$
与 $X_2$ 有 $X_3$
北韩    邦交

$X_1$
have $X_3$ with $X_2$
diplomatic    North
relations    Korea

# Extracting Hiero rules



X → 与 北 韩 有 邦交, have diplomatic relations with North Korea

X → 与 北 韩 有 邦交,
have diplomatic relations
with North Korea

X → 邦交,
   diplomatic relations

# Extracting Hiero rules



X → 与 北 韩 有 邦 交,
have diplomatic relations
with North Korea

X → 邦 交,
diplomatic relations

X → 北 韩,
North Korea

15

X → 与 北 韩 有 邦交,
have diplomatic relations
with North Korea

X → 邦交,
    diplomatic relations

X → 北 韩,
    North Korea

X → 与 北 韩 有 邦交,
have diplomatic relations with North Korea

X → 邦交,
  diplomatic relations

X → 北 韩,
  North Korea

X → 与 $X_1$ 有 $X_2$,
  have $X_2$ with $X_1$

15

# Extracting Hiero rules



X → 与 北 韩 有 邦交,
have diplomatic relations with North Korea

X → 邦交,
  diplomatic relations

X → 北 韩,
  North Korea

X → 与 X₁ 有 X₂,
  have X₂ with X₁

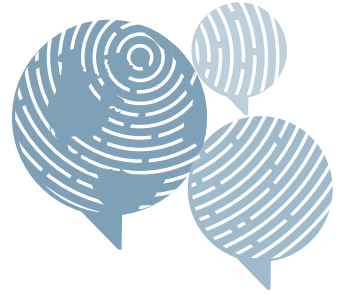VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea
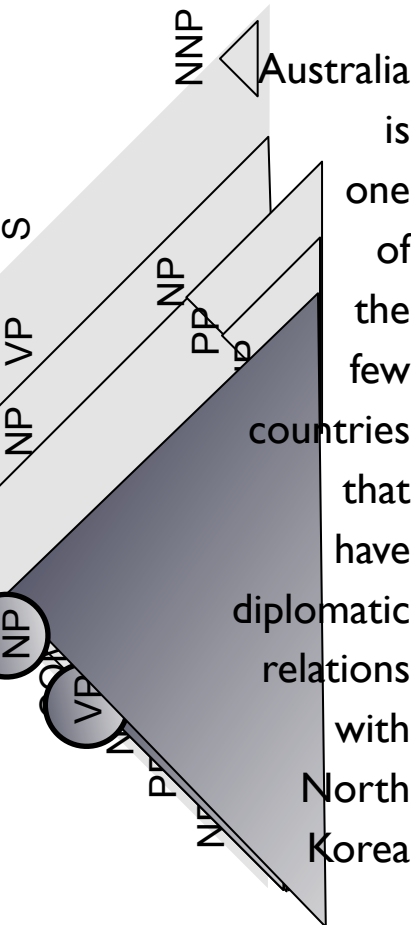
VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea

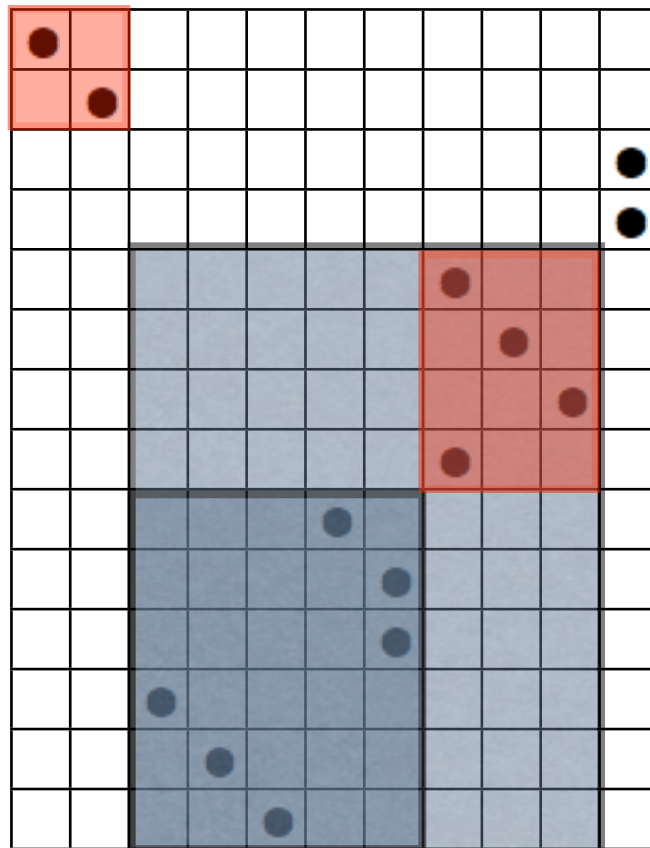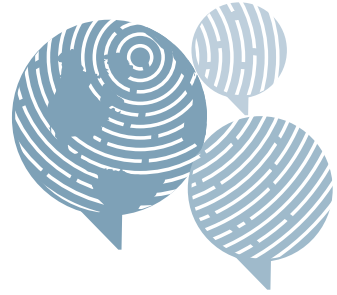??? → 的 少数 国家, the few countries that

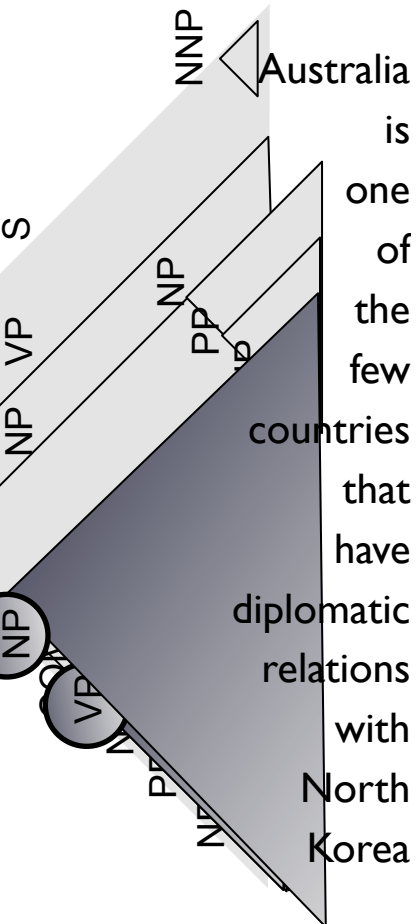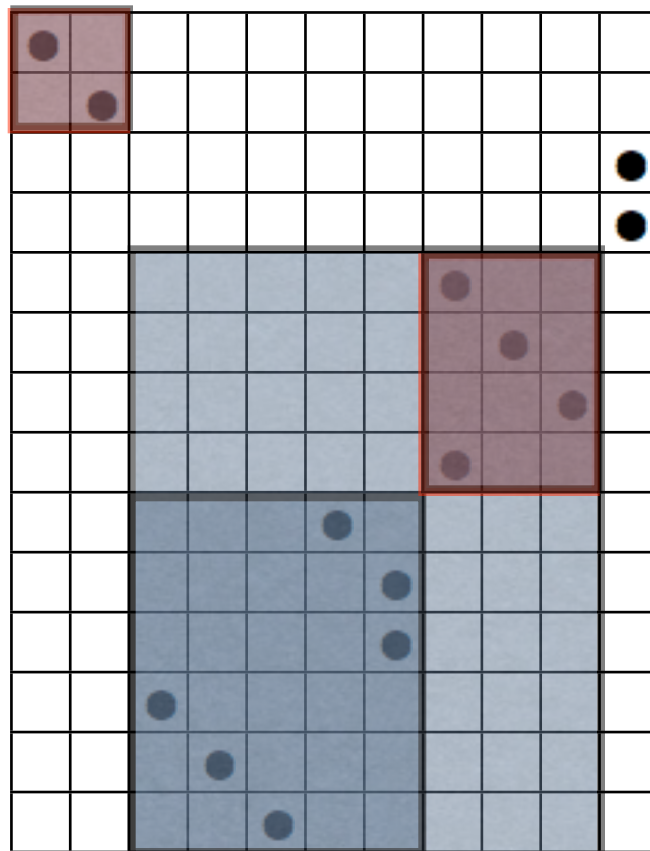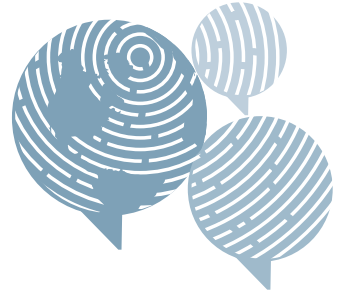VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea

??? → 的 少数 国家, the few countries that

??? → 澳洲 是, Australia is

VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

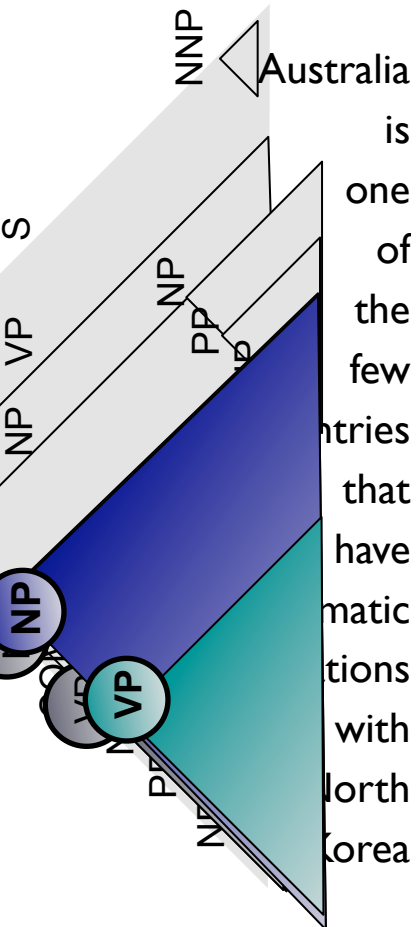NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea

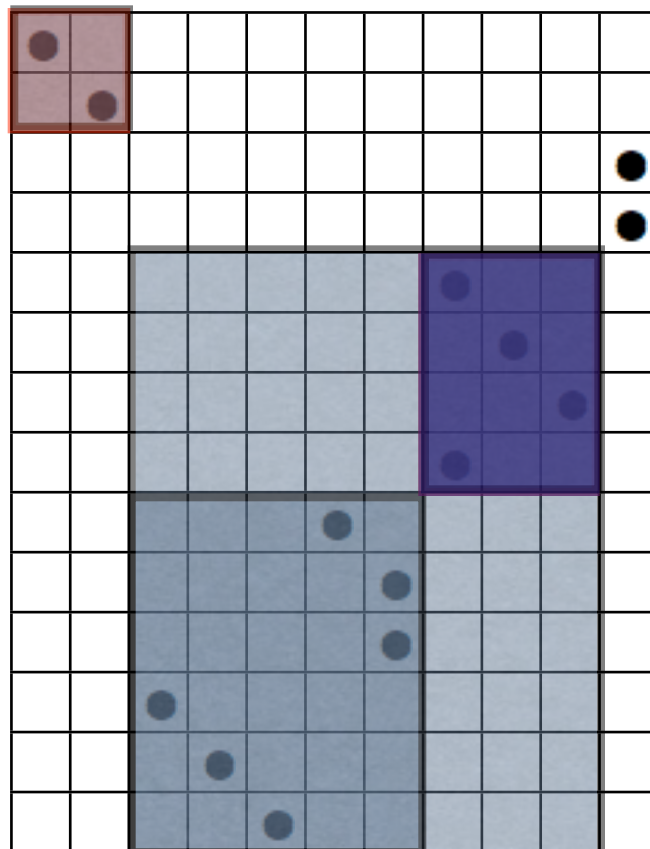??? → 的 少数 国家, the few countries that

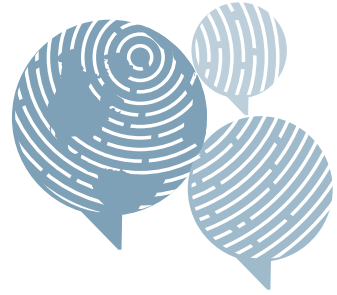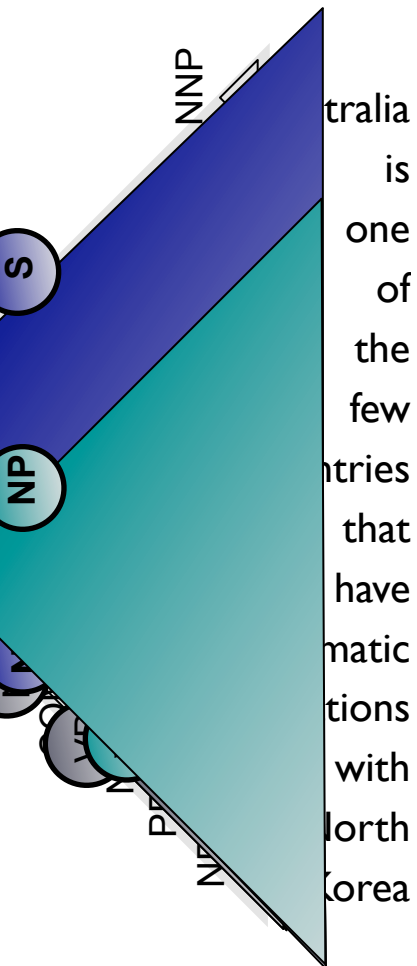??? → 澳洲 是, Australia is

# Extracting Syntactic Rules



VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea
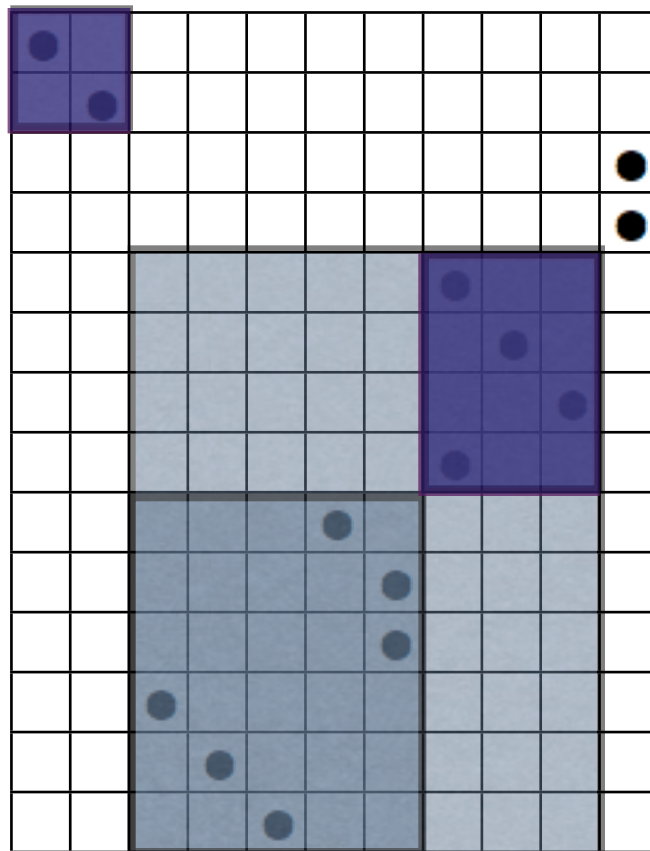
NP/VP → 的 少数 国家, the few countries that

??? → 澳洲 是, Australia is

# Extracting Syntactic Rules

VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea

NP/ VP → 的 少数 国家, the few countries that
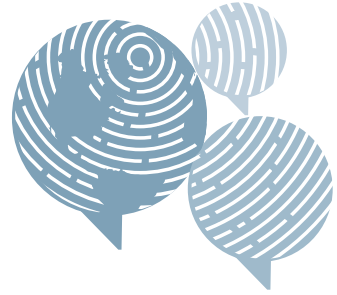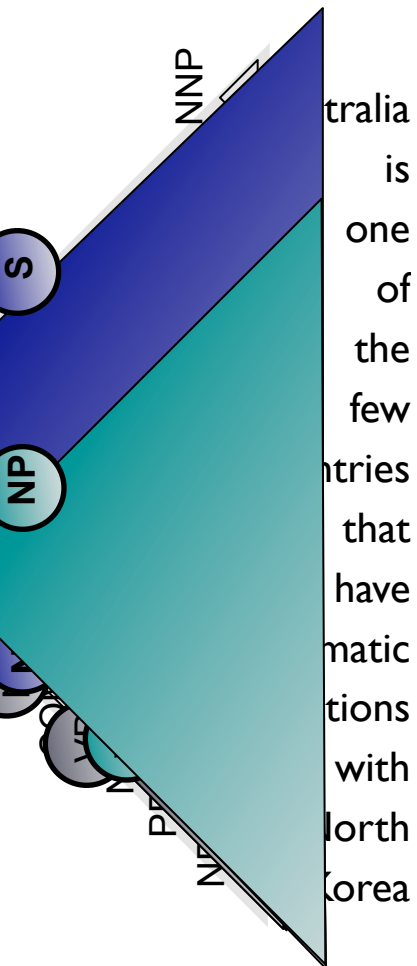
S/ NP → 澳洲 是, Australia is

# Extracting Syntactic Rules



VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

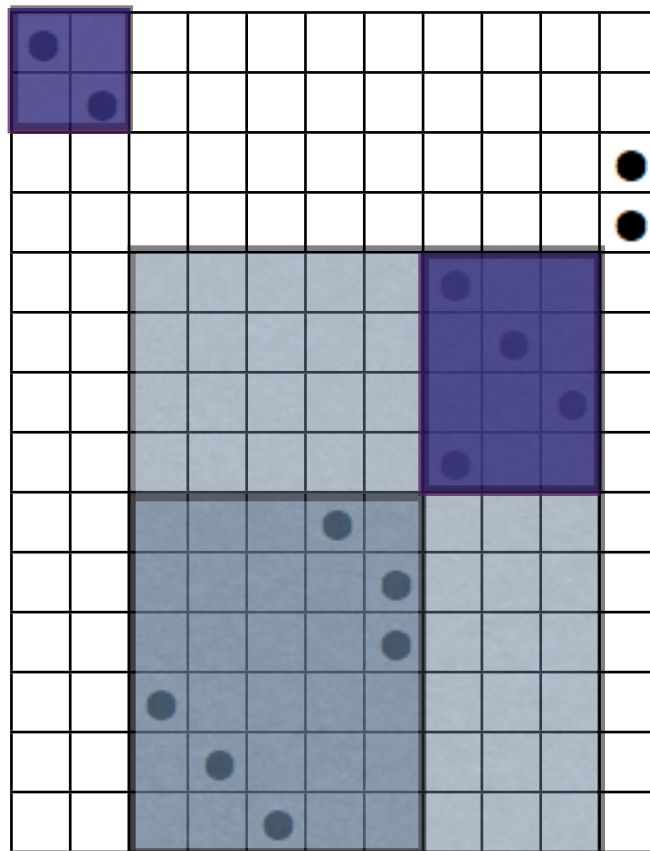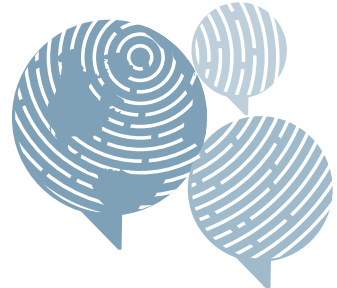NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea

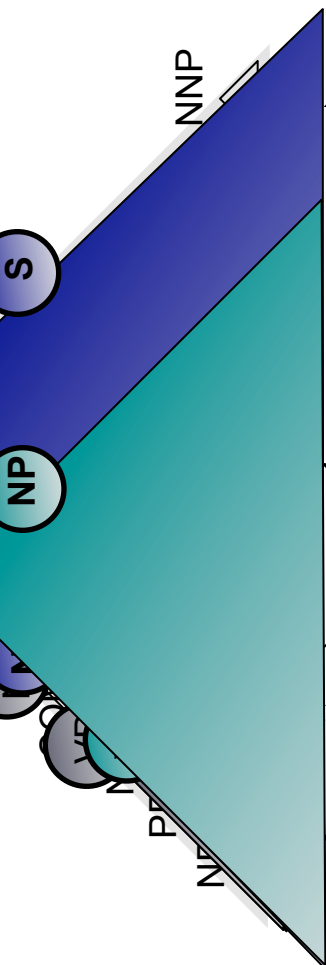NP/VP → 的 少数 国家, the few countries that

S/NP → 澳洲 是, Australia is

VP → 与 北 韩 有 邦交, have diplomatic relations with North Korea

NP → 与 北 韩 有 邦交 的 少数 国家, the few countries that have diplomatic relations with North Korea

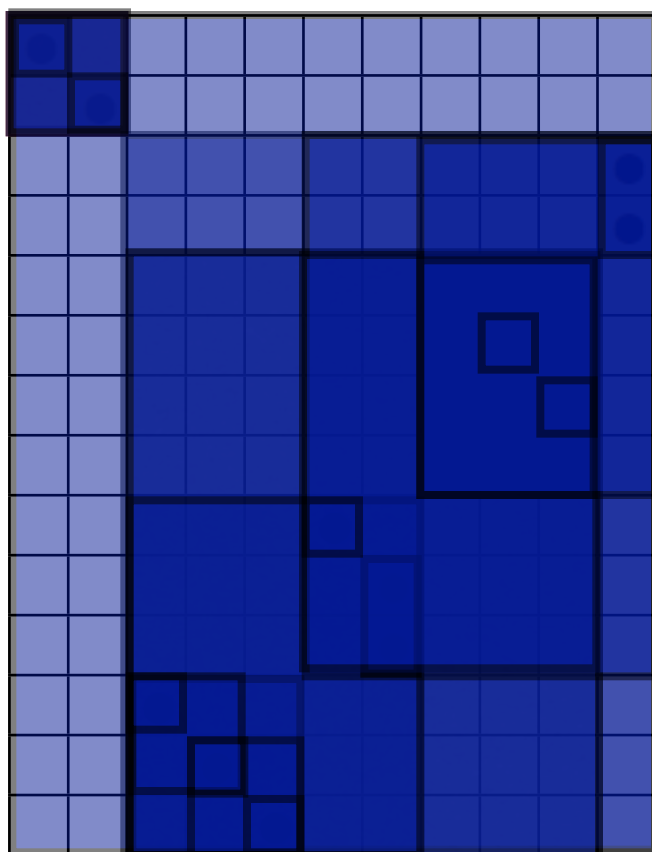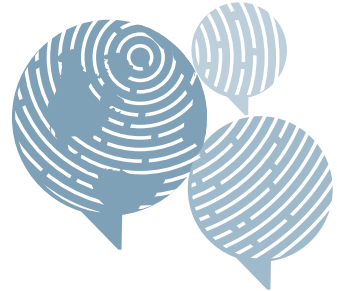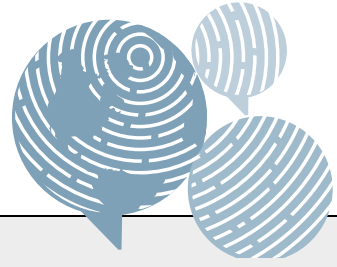NP/VP → 的 少数 国家, the few countries that
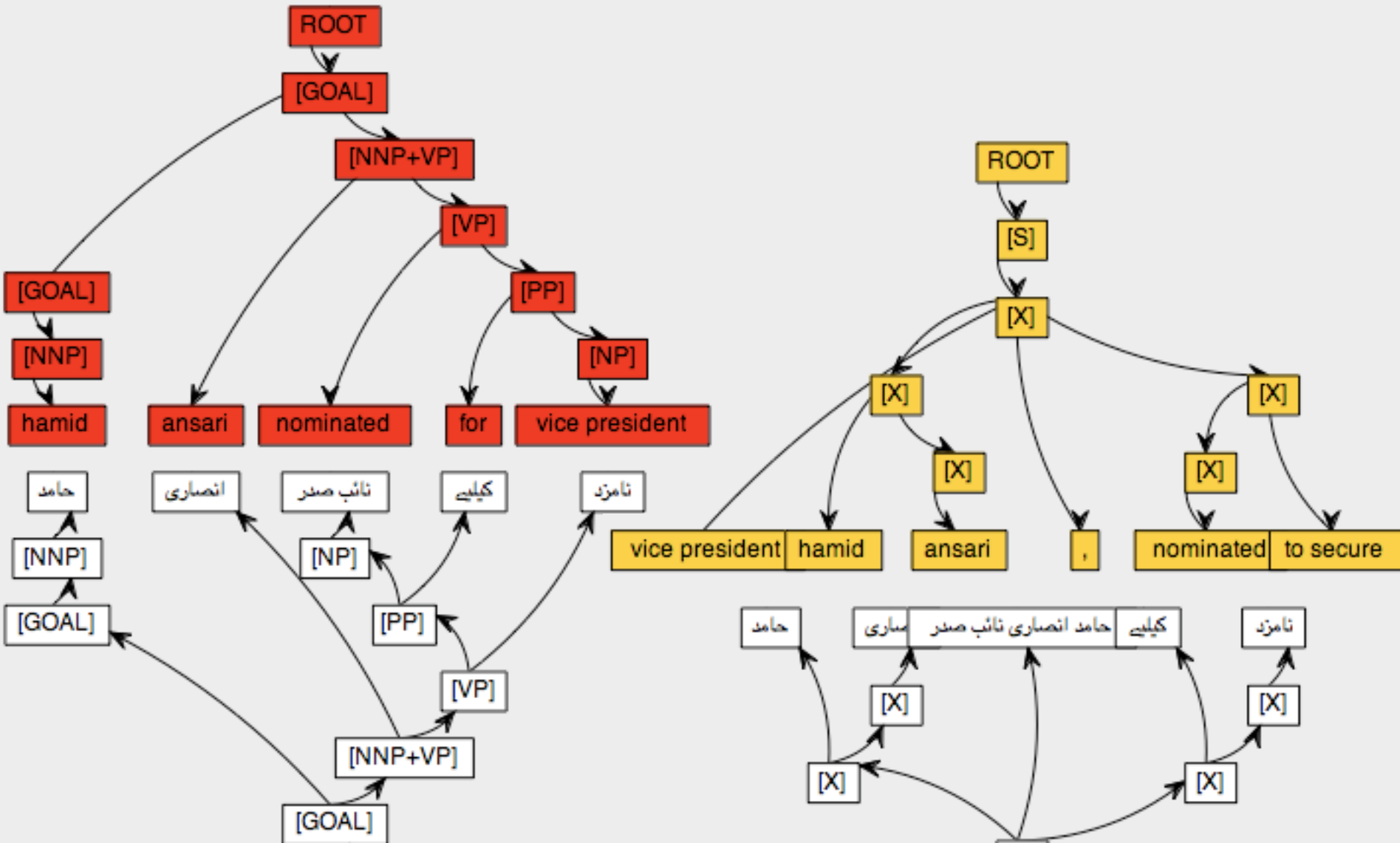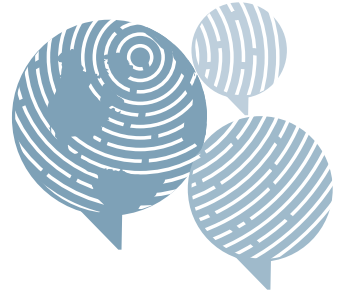
S/NP → 澳洲 是, Australia is

# New training paradigm

- Training data: word-aligned bilingual parallel corpus, with parse trees
  - No need to parse the Urdu, just parse the English
  - Method is therefore transferable to other resource poor languages
- Extract SCFG rules with syntactic nonterminals
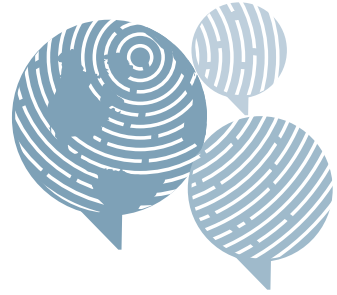- For non-constituent phrases use CCG-style nonterminals
- Same coverage as Hiero model

ROOT

[GOAL]

[NNP+VP]

[VP]

[PP]

[GOAL]

[NNP]

[NP]

hamid | ansari | nominated | for | vice president

حامد | انصاری | نائب صدر | کیلیے | نامزد

[NNP] | [NP]

[GOAL] | [PP]

[VP]

[NNP+VP]

[GOAL]

---

ROOT

[S]

[X]

[X] | [X]

[X] | [X]

vice president | hamid | ansari | , | nominated | to secure

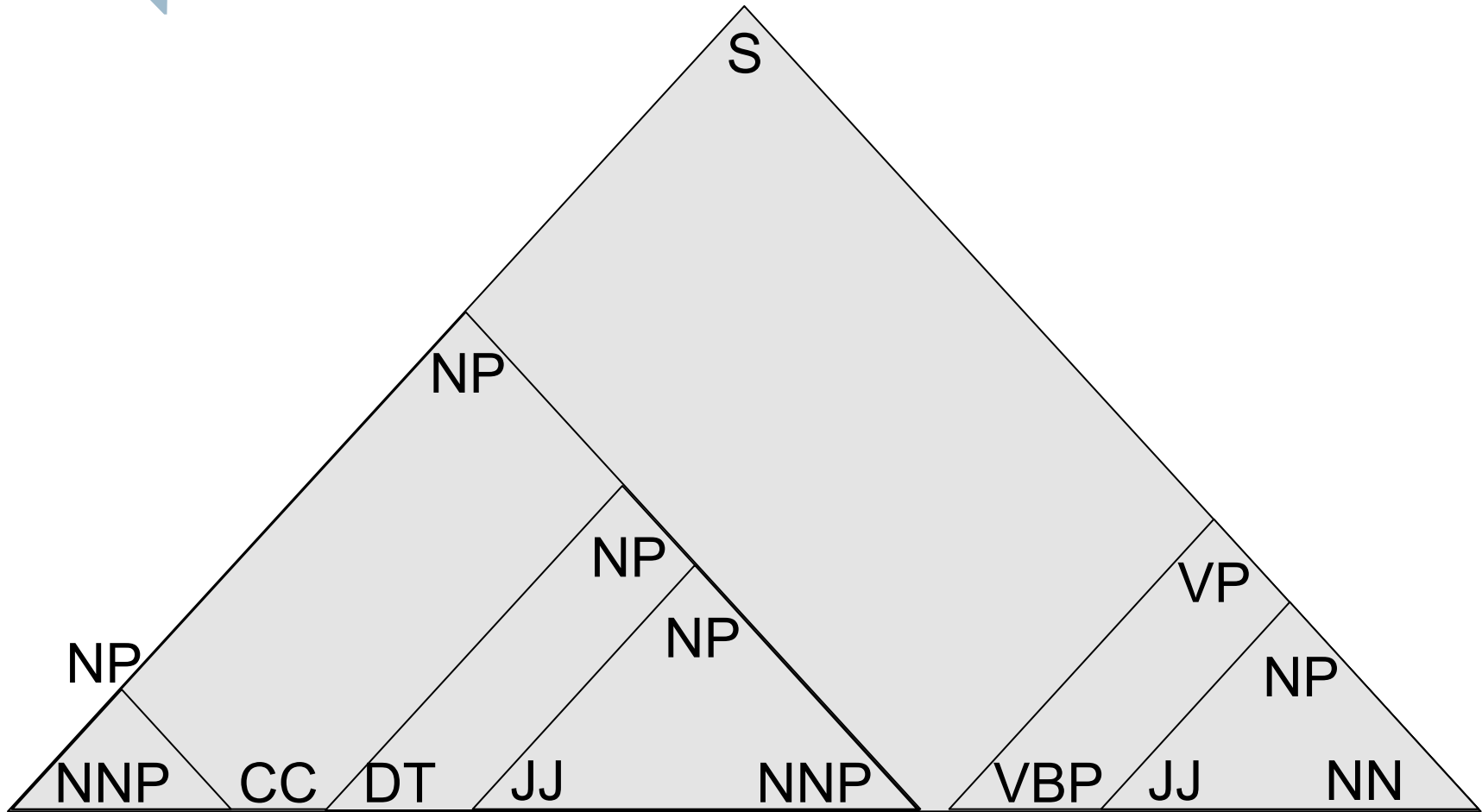حامد | کیلیے | نائب صدر انصاری حامد | نامزد
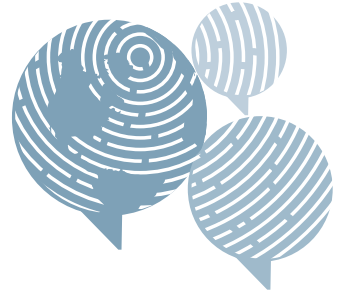
[X] | [X]

[X] | [X]

# Integrating HIVEs into SCFGs

- A primary goal of the workshop was to include semantic entities in the translation process

- Semantic entities included high information value elements such as

  - Named entities

  - Modality

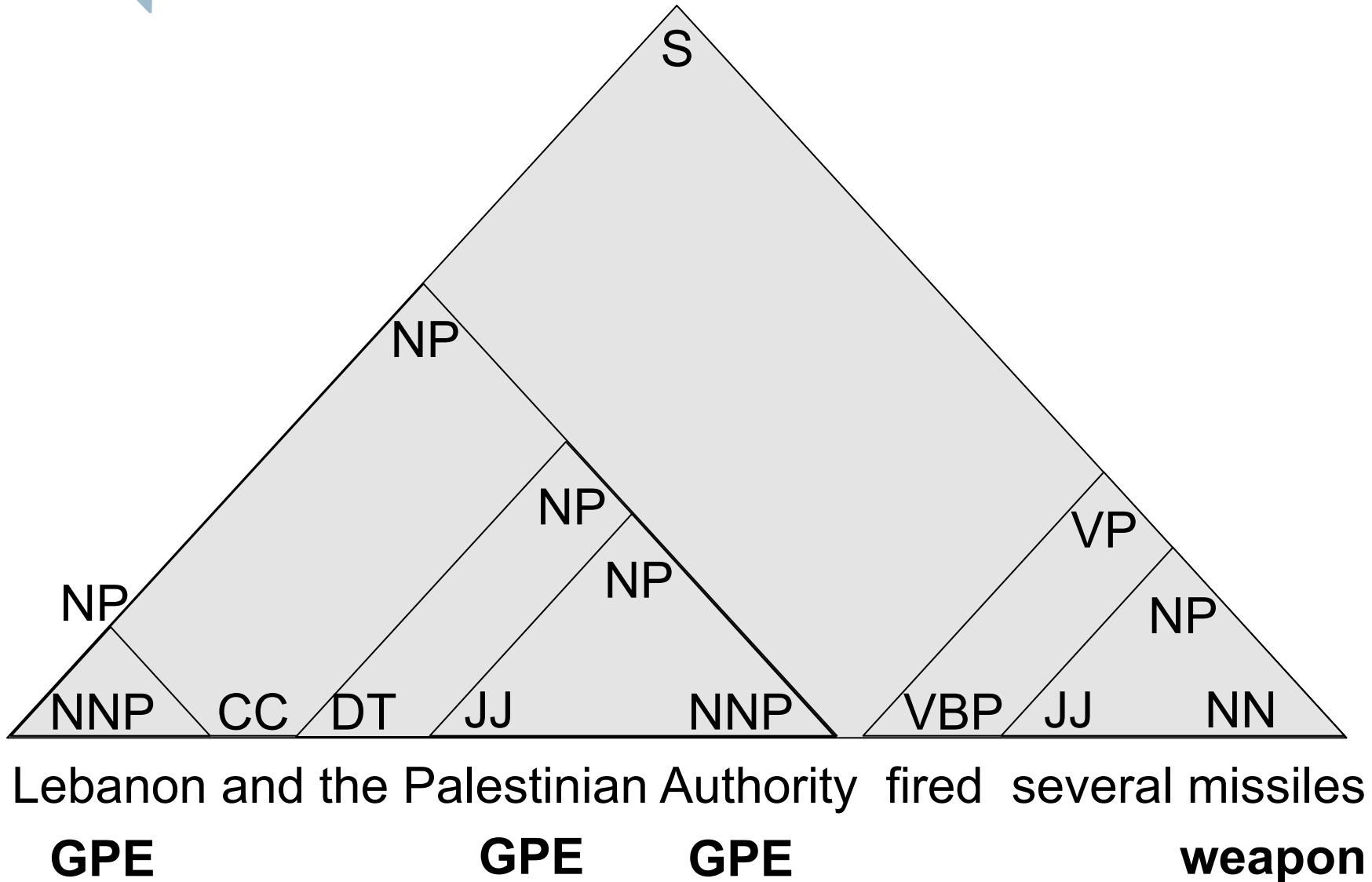- These elements are included in SCFGs by grafting onto parse trees

S

NP

NP

VP

NP

NP

NP

NNP    CC    DT    JJ    NNP    VBP    JJ    NN
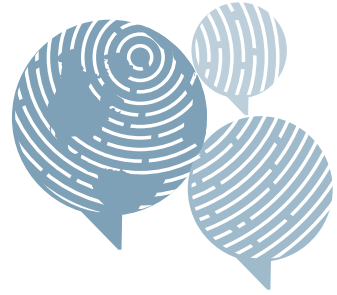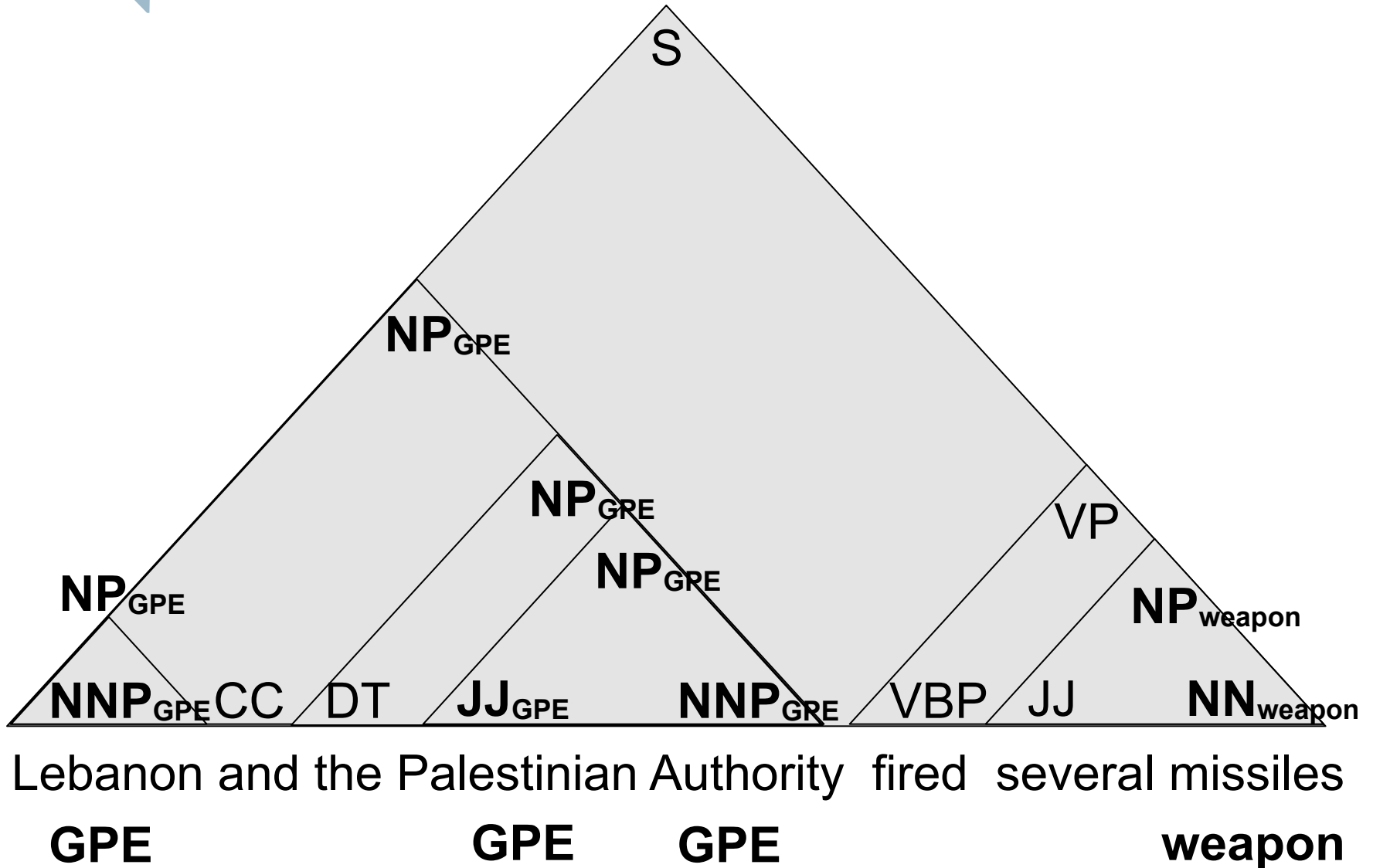
Lebanon and the Palestinian Authority  fired  several missiles

# Example tree graft

# Example tree graft

# Simple Urdu modality tagger

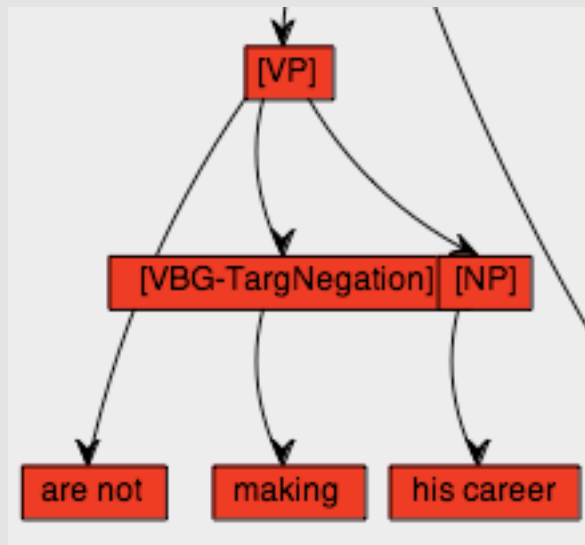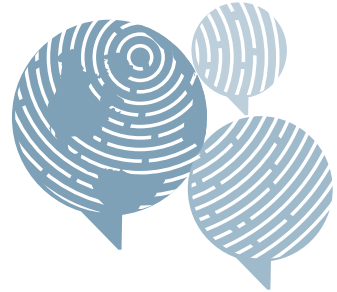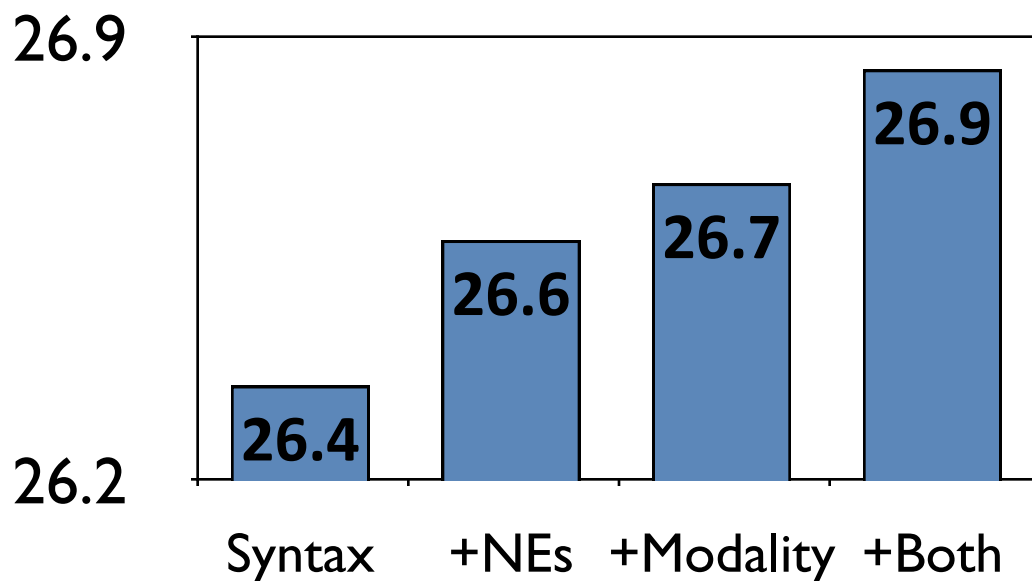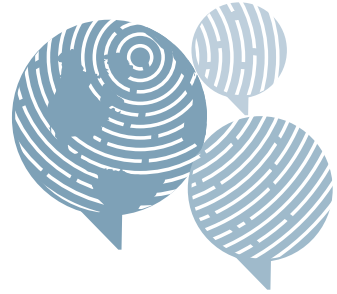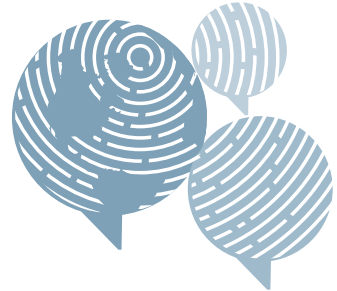# HIVE improvements

- HIVEs improve quality over using syntax alone by 0.5 Blue points

- NEs better capture distributional properties of people, places, organizations (over generic "NP").  Modality better structures VPs.

# Conclusions

- Using syntax-based translation models resulted in huge improvements in quality

- Urdu is an ideal language to show off the advantages of syntax
  - Very small amount of training data
  - Very different word order than English

- No need for any specialized Urdu tools
  - Therefore applicable to other languages
  - Likely to find similar gains for low-resource languages that have different word order

# Free software!

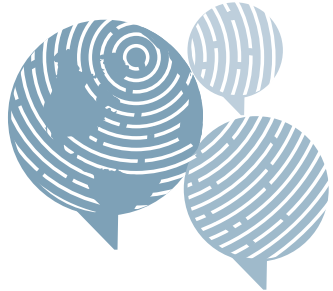- The Joshua Decoder is open source -- try it out today! http://cs.jhu.edu/~ccb/joshua/

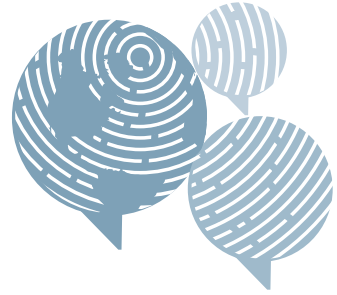- Zhifei Li, one of our lead developers, is graduating -- hire him!

# Thanks!

# Breakdown of improvements



Bar chart values:
- Hiero trigram l...: 20.9
- Hiero w/giza: 20.9
- Moses Giza + 5...: 22.3
- Hiero 5-gram: 23.1
- Syntax w/3 FFs: 18.0
- Syntax w/8 FFs: 20.0
- Syntax w/12 FF...: 26.6
- + NEs: 27.0
- + multiple align...: 27.3
- + Min Bayes Ris...: 27.5
- + strip UNKs: 28

# NIST Eval Results

X → 与 北 韩 有 邦交, have diplomati
with North Ko

澳 是 与 北 韩 有 邦 的 少 国 之
洲　　　　　　交　数 家 一

| | 澳洲 | 是 | 与 | 北 | 韩 | 有 | 邦交 | 的 | 少数 | 国家 | 之一 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Australia | ● | | | | | | | | | | |
| is | | ● | | | | | | | | | |
| one | | | | | | | | | | | ● |
| of | | | | | | | | | | | ● |
| the | | | | | | | ● | | | | |
| few | | | | | | | | ● | | | |
| countries | | | | | | | | | | ● | |
| that | | | | | | | ● | | | | |
| have | | | | | | ● | | | | | |
| diplomatic | | | | | | | ● | | | | |
| relations | | | | | | | ● | | | | |
| with | | | ● | | | | | | | | |
| North | | | | ● | | | | | | | |
| Korea | | | | | ● | | | | | | |

X → 与 北 韩 有 邦交, have diplomati
with North Ko

X → 邦交, diplomatic relations

# Extracting Hiero rules



X → 与 北 韩 有 邦交, have diplomati[...]
with North Ko[...]

X → 邦交, diplomatic relations

X → 北 韩, North Korea

X → 与 北 韩 有 邦交, have diplomati[c]
 with North Ko[rea]

X → 邦交, diplomatic relations

X → 北 韩, North Korea

X → 与 北 韩 有 邦交, have diplomati
with North Ko

X → 邦交, diplomatic relations

X → 北 韩, North Korea

X → 与 $X_1$ 有 $X_2$, have $X_2$ with $X_1$

X → 与 北 韩 有 邦交, have diplomatic
with North Ko

X → 邦交, diplomatic relations

X → 北 韩, North Korea

X → 与 X₁ 有 X₂, have X₂ with X₁

澳 是 与 北 韩 有 邦 的 少 国 之
洲　　　　　　交　数 家 一

Australia
is
one
of
the
few
countries
that
have
diplomatic
relations
with
North
Korea

NNP
S
VP
NP
PP
NP
NP
COMP
VP
NP

VP

VP → 与 北 韩 有 邦交,
　　　have diplomatic relation

NP → 与 北 韩 有 邦交 的 少
　　　the few countries that
　　　relations with North K

NP → VP 的 少数 国家,
　　　the few countries that

S/NP → 澳洲, Australia is

S/NP → 澳洲，Australia is

# Example tree graft



NNP  ,  NP VBZ RB  DT VBN      NN    TO NNP  NNP  , VBD

Nance , who  is  also  a  paid consultant to ABC News , said ....
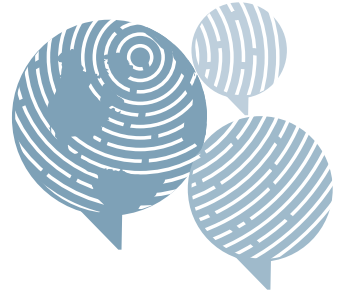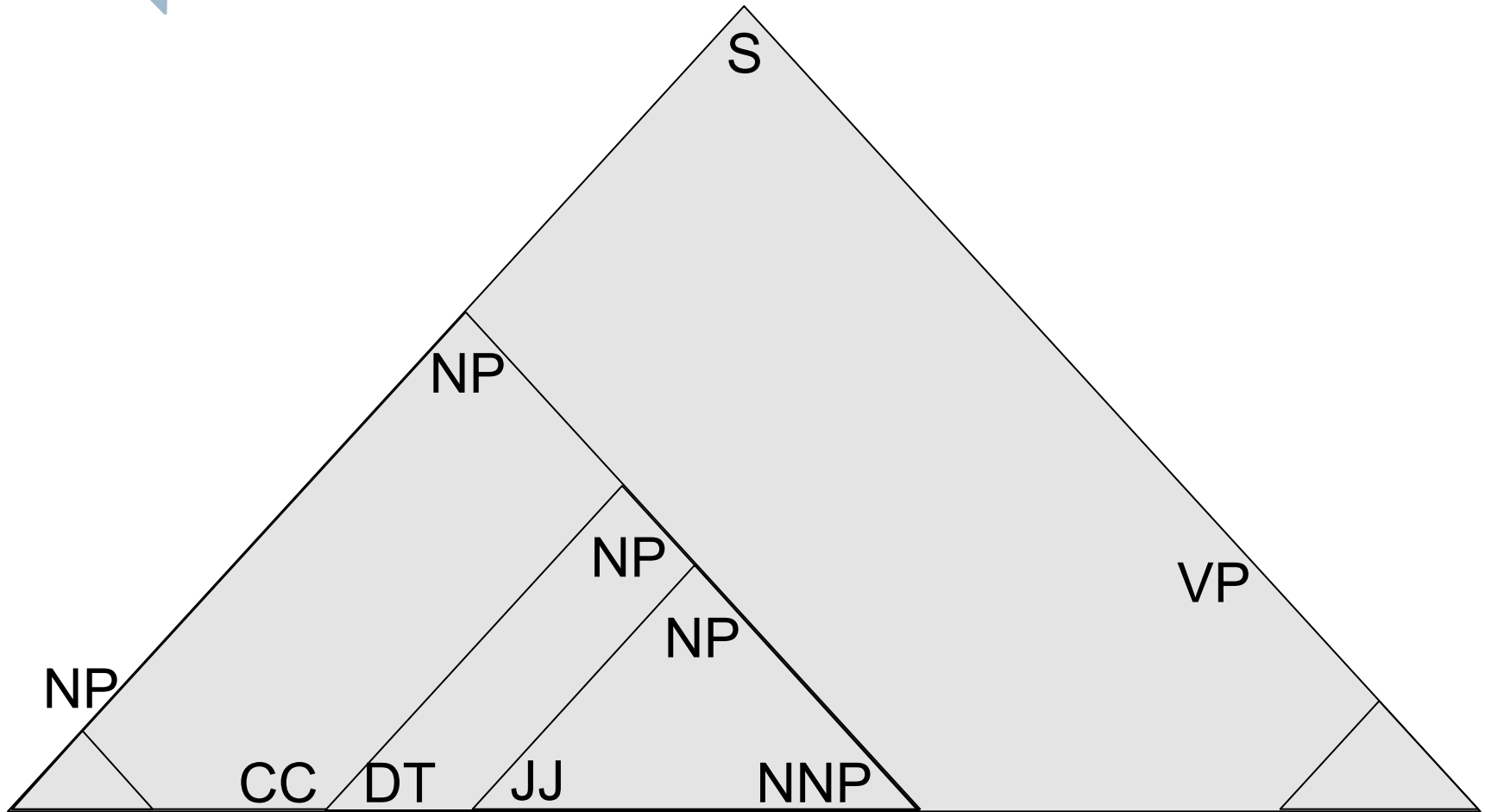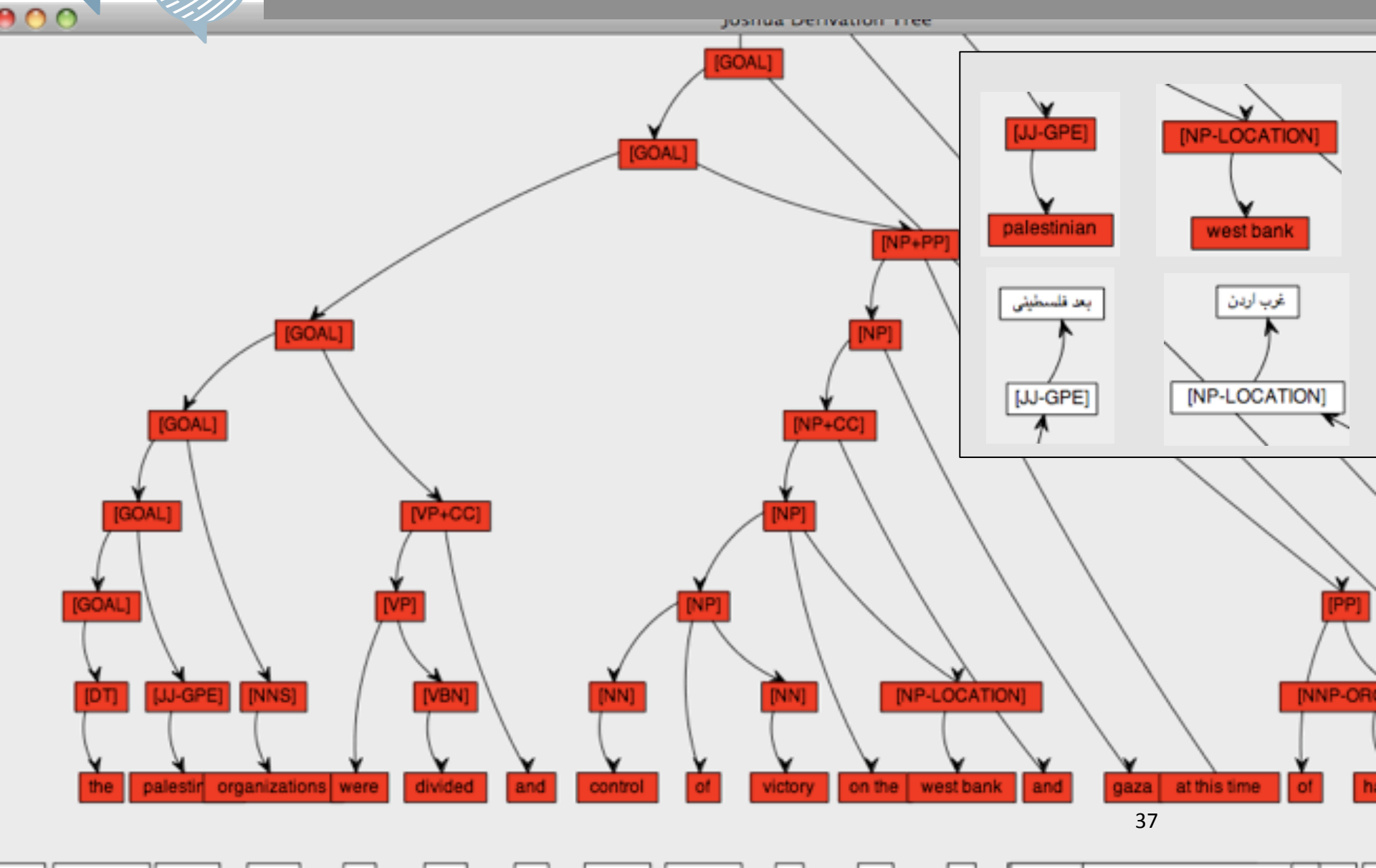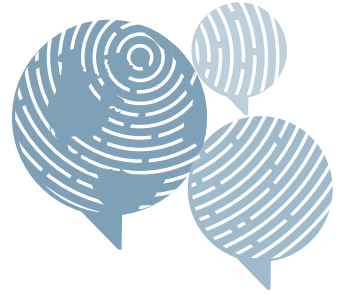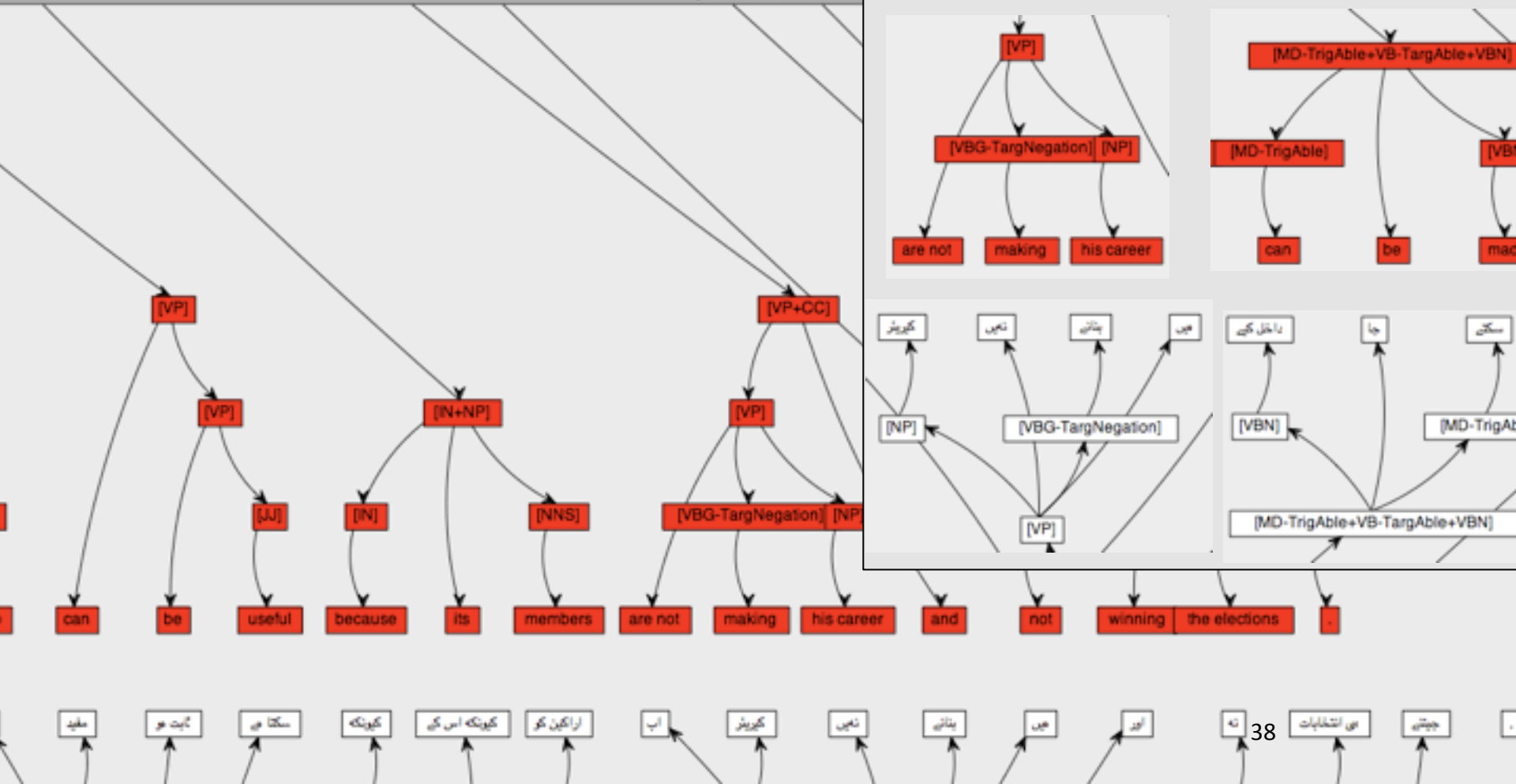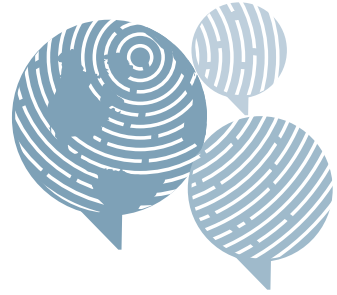
## Baseline

### 'first nuclear experiment in 1990 was'

Thomas red Unilever National Laboratory of the United States in ویپن designer, are already working on the book of Los ایلموس National Laboratory ڈینی, former director of the technical انٹیلجنس written with the cooperation of سٹلمین.

This book 'nuclear express: political history and the expansion of bomb' has been written, and the two writers have also claimed that the country has made nuclear bomb is he or any other country's nuclear secrets to چرائے or that of any other nuclear power cooperation is achieved.

Thomas Reid said in a news یوایس interview that in 1990 in the era of Benazir Bhutto China had the experience of Pakistan's first nuclear bomb.
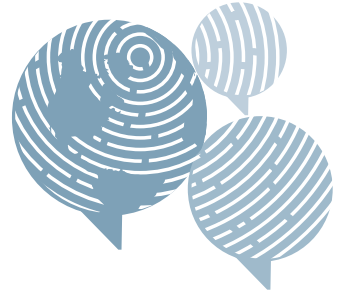
### The First Nuclear T

Thomas red of the U Laboratory in design the book of Los Alan former director of the the cooperation of D

This book under the expressway: the poli this has been writte claimed that the cou bomb or any other c secrets, or any of the cooperation.

Thomas Reid said in that Benazir Bhutto i

## Source

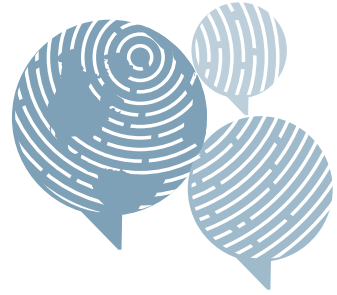| | |
|---|---|
| ناگاؤں نے آسام میں آگ لگا دی<br><br>بدھ کے روز مشتعل ناگا قبائلیوں نے منی پور کے دس سکولوں کو بھی نذر آتش کردیا تھا۔<br><br>پولیس کے مطابق سینکڑوں کی تعداد میں ناگالینڈ کے مسلح قبائلیوں نے آسام کے گلیکی اور سیبسا گر کے تین گاؤں میں آگ لگا دی۔<br><br>اس حملہ کے بعد بڑی تعداد میں مقامی باشندوں نے علاقوں کو خالی کردیا ہے۔<br><br>ناگالینڈ دعویٰ کرتا ہے کہ ریاست آسام اس کے بعض خطوں پر قابض ہے۔ | **Nagas Set Fire** i<br><br>On Wednesday, schools in Manip<br><br>According to poli tribesmen of Nag Gulleki and Sisag<br><br>A large number o area after this att<br><br>Nagaland claims some of its territo<br><br>While Assam sta |

Reference          pre-SCALE

**Nagas Set Fire in Assam**          **Has Imposed a Fire in**

On Wednesday, angry Naga tribesmen set 10 schools in Manipur on fire.

On Wednesday, the triba schools was also burnt.

According to police, hundreds of armed tribesmen of Nagaland set three villages of Gulleki and Sisagar in Assam.

According to the police, hundreds of armed tribe assam and three set the

A large number of natives have vacated the area after this attack.

After this attack. Local r numbers to the areas.

Nagaland claims that Assam state is occupying some of its territory.

Claim of assam. That th

While Assam state says that Nagaland has
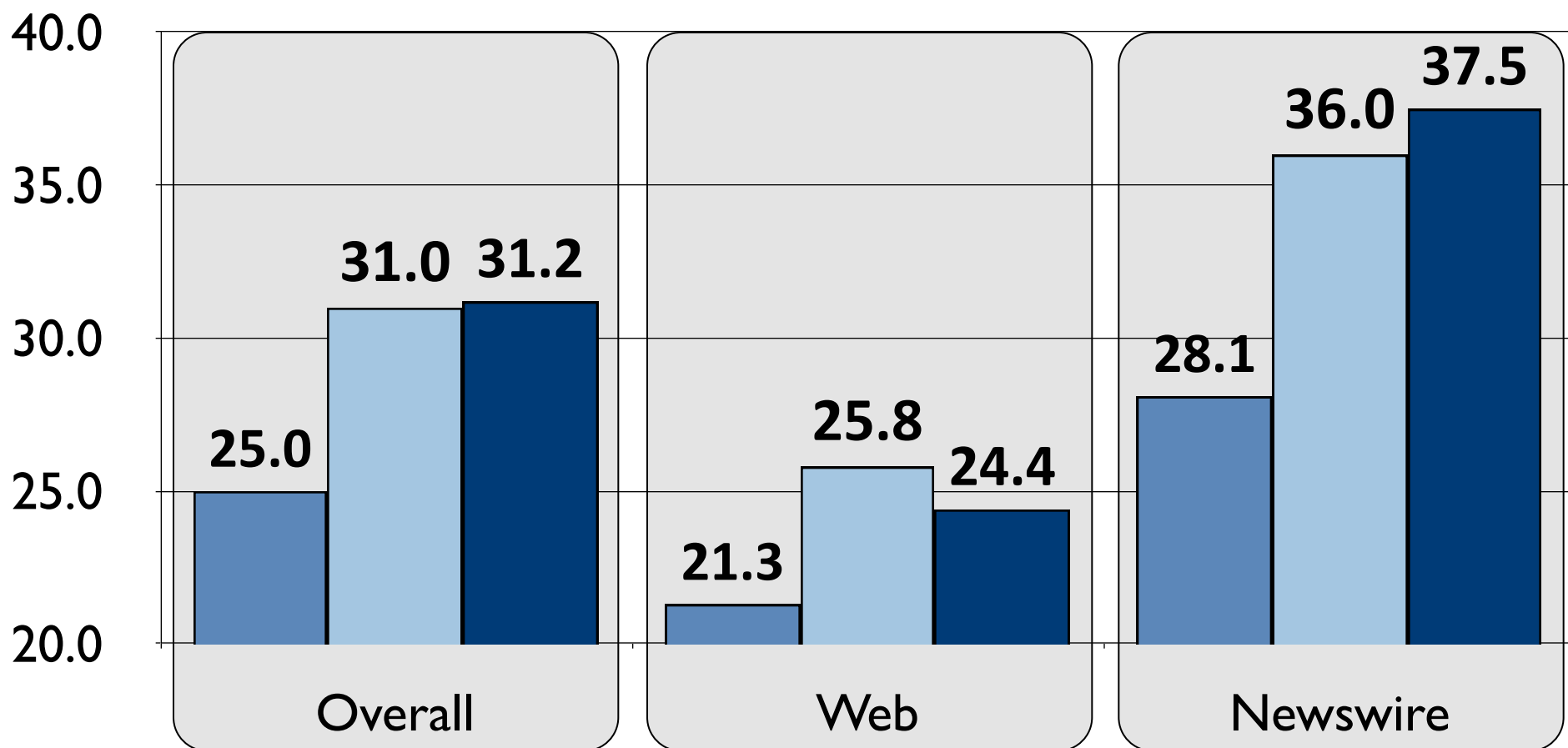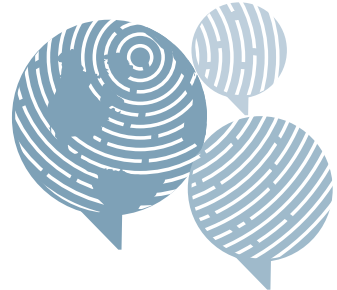
While of assam has said areas of his into custody

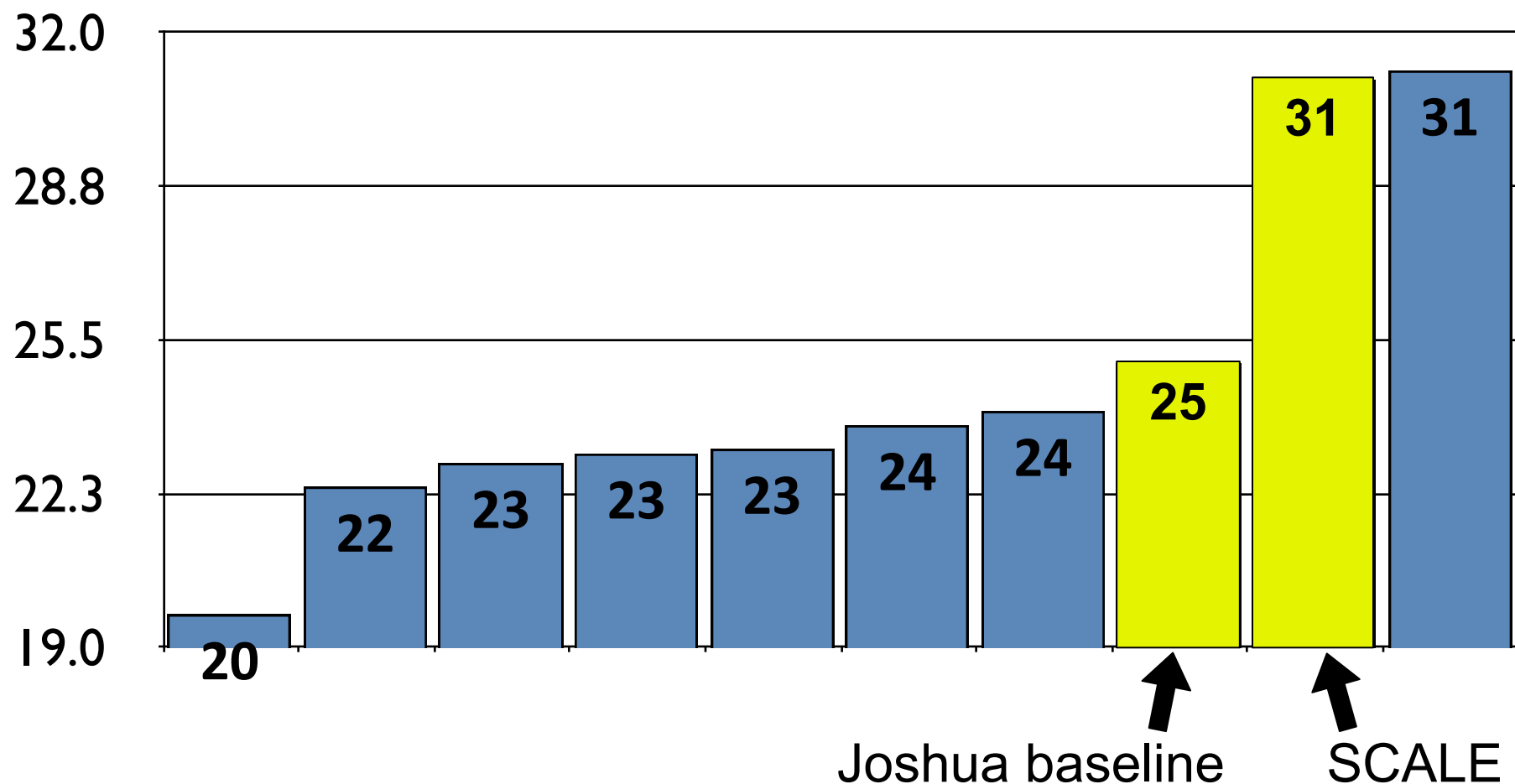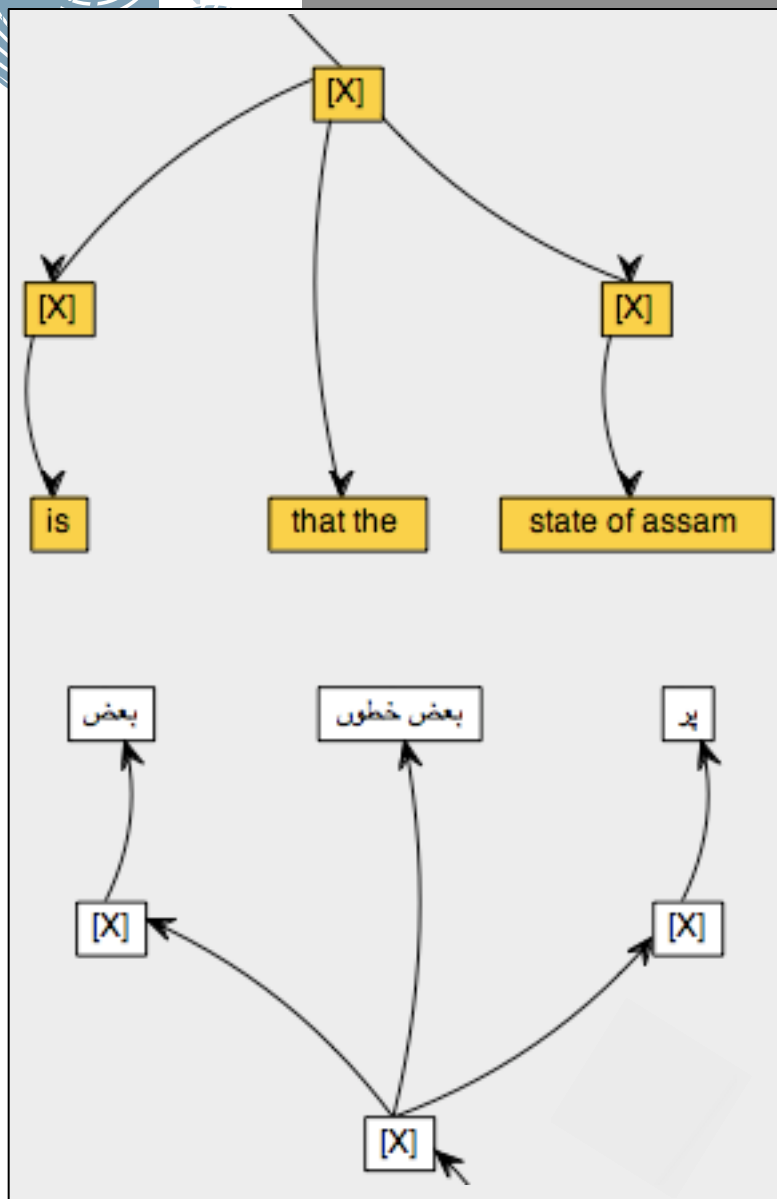# Breakdown of improvements

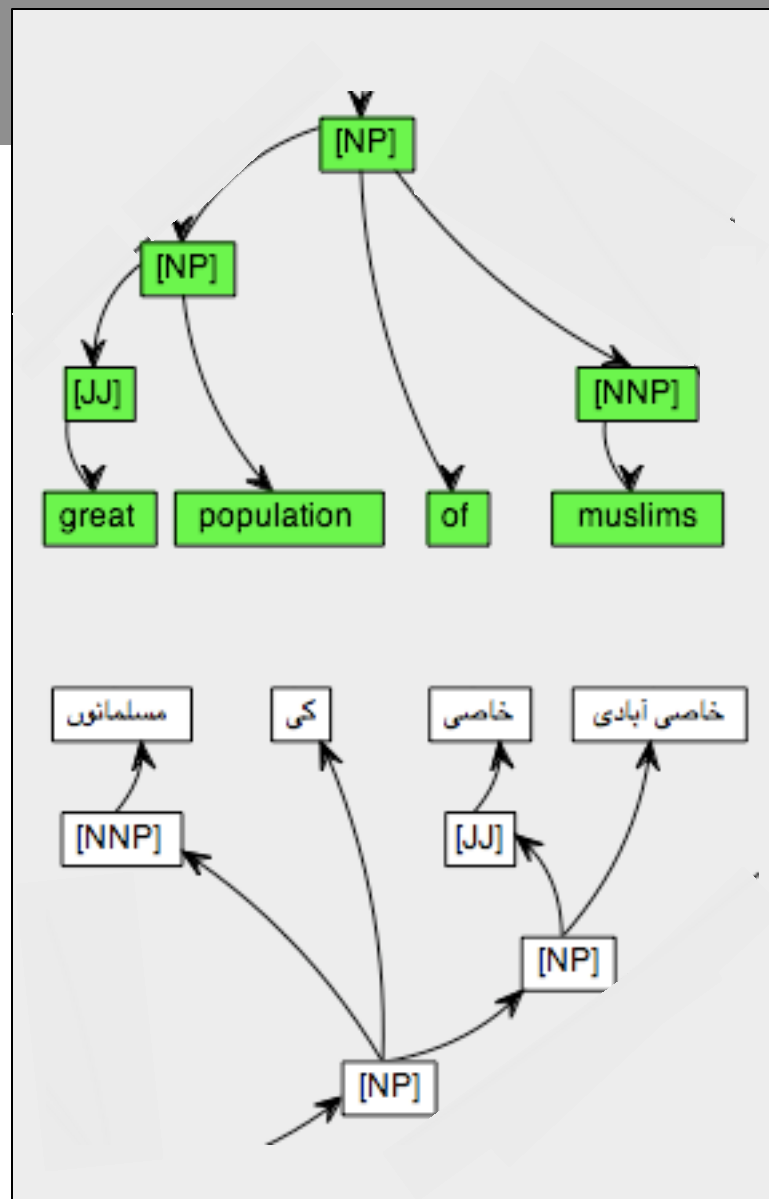# NIST Eval Results
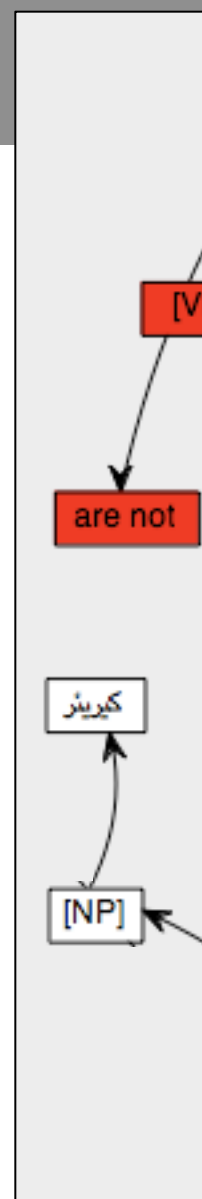
# The field for the Urdu task

Standard Hierarchical Rules          Syntactic Enhancements          Sema