

Statistical Machine Translation of Natural Languages

Heiko Vogler
Technische Universität Dresden
Germany

Graduiertenkolleg “Quantitative Logics and Automata”
Dresden, November, 2012

Weighted Tree Automata and Weighted Tree Transducers

can help in

Statistical Machine Translation of Natural Languages

Heiko Vogler

Technische Universität Dresden

Germany

Graduiertenkolleg “Quantitative Logics and Automata”
Dresden, November, 2012

outline of the talk:

- ▶ Statistical machine translation

outline of the talk:

- ▶ Statistical machine translation
- ▶ Modeling with wta and wtt

outline of the talk:

- ▶ Statistical machine translation
- ▶ Modeling with wta and wtt
- ▶ Using automata theoretic results to “improve” modeling

outline of the talk:

- ▶ Statistical machine translation
- ▶ Modeling with wta and wtt
- ▶ Using automata theoretic results to “improve” modeling
- ▶ Summary

outline of the talk:

- ▶ **Statistical machine translation** **no survey!**
- ▶ Modeling with wta and wtt
- ▶ Using automata theoretic results to “improve” modeling
- ▶ Summary

given:

- ▶ source language SL
- ▶ target language TL

find:

translation $h : SL \rightarrow TL$

e.g.

SL = English

TL = German

s = I saw the man with the telescope

$h(s)$ = Ich sah den Mann durch das Tel.

given:

- ▶ source language SL
- ▶ target language TL

find:

machine translation $h : SL \rightarrow TL$

e.g.

SL = English

TL = German

s = I saw the man with the telescope

$h(s)$ = Ich sah den Mann durch das Tel.

given:

- ▶ source language SL
- ▶ target language TL

find:

machine translation $h : SL \rightarrow TL$

e.g.

SL = English

TL = German

s = I saw the man with the telescope

$h(s)$ = Ich sah den Mann durch das Tel.

machine translation \rightsquigarrow statistical machine translation

given:

- ▶ source language SL
- ▶ target language TL

find:

machine translation $h : SL \rightarrow TL$

e.g.

SL = English

TL = German

s = I saw the man with the telescope

$h(s)$ = Ich sah den Mann durch das Tel.

machine translation \rightsquigarrow statistical machine translation

modeling - training - evaluation

given:

- ▶ source language SL
- ▶ target language TL

find:

machine translation $h : SL \rightarrow TL$

e.g.

SL = English

TL = German

s = I saw the man with the telescope

$h(s)$ = Ich sah den Mann durch das Tel.

machine translation \rightsquigarrow statistical machine translation

modeling - training - evaluation

given:

- ▶ source language SL
- ▶ target language TL

find:

machine translation $h : SL \rightarrow TL$

e.g.

SL = English

TL = German

$s =$ I saw the man with the telescope

$h(s) =$ Ich sah den Mann durch das Tel.

machine translation \rightsquigarrow **statistical** machine translation

modeling - training - evaluation

assumptions \longrightarrow modeling \longrightarrow \mathcal{H} hypothesis space

assumptions: mental work, experience, no data

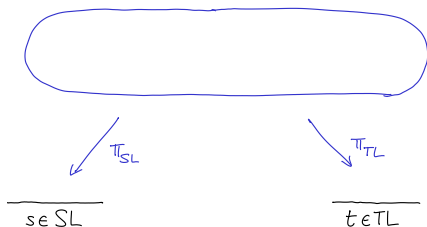
hypothesis space: $\mathcal{H} \subseteq \{h \mid h : SL \rightarrow TL\}$

log-linear modeling:

$$\frac{\quad}{seSL}$$

$$\frac{\quad}{teTL}$$

log-linear modeling:

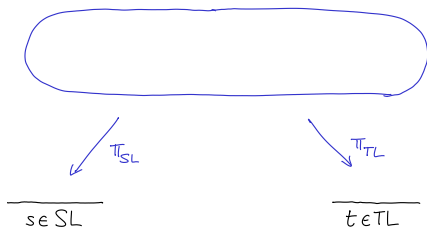


set Y of correspondence
structures [Liang et al. 06]

$$\pi_{SL} : Y \rightarrow SL$$

$$\pi_{TL} : Y \rightarrow TL$$

log-linear modeling:



set Y of correspondence structures [Liang et al. 06]

$$\pi_{SL} : Y \rightarrow SL$$

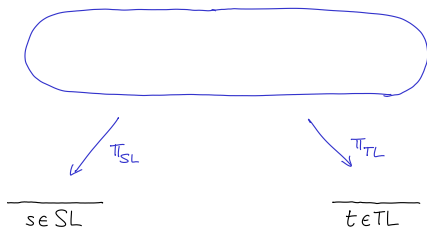
$$\pi_{TL} : Y \rightarrow TL$$

hypothesis space: $\mathcal{H} = \{h_{\lambda, \Phi} \mid \lambda \in \mathbb{R}_{\geq 0}^m, \Phi : Y \rightarrow \mathbb{R}^m\}$

$$h_{\lambda, \Phi} : SL \rightarrow TL$$

$$s \mapsto \pi_{TL} \left(\operatorname{argmax}_{y \in Y: \pi_{SL}(y)=s} \lambda \cdot \Phi(y) \right)$$

log-linear modeling:



set Y of correspondence structures [Liang et al. 06]

$$\pi_{SL} : Y \rightarrow SL$$

$$\pi_{TL} : Y \rightarrow TL$$

hypothesis space: $\mathcal{H} = \{h_{\lambda, \Phi} \mid \lambda \in \mathbb{R}_{\geq 0}^m, \Phi : Y \rightarrow \mathbb{R}^m\}$

$$h_{\lambda, \Phi} : SL \rightarrow TL$$

$$s \mapsto \pi_{TL} \left(\operatorname{argmax}_{y \in Y: \pi_{SL}(y)=s} \lambda \cdot \Phi(y) \right)$$

$$\lambda_1 \cdot \Phi(y)_1 + \dots + \lambda_m \cdot \Phi(y)_m$$

hypothesis space: $\mathcal{H} = \{h_{\lambda, \Phi} \mid \lambda \in \mathbb{R}_{\geq 0}^m, \Phi : Y \rightarrow \mathbb{R}^m\}$

$$h_{\lambda, \Phi} : \text{SL} \rightarrow \text{TL}$$

$$s \mapsto \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(y)=s}{\operatorname{argmax}}_{y \in Y} \lambda \cdot \Phi(y) \right)$$

$$\lambda_1 \cdot \Phi(y)_1 + \dots + \lambda_m \cdot \Phi(y)_m$$

hypothesis space: $\mathcal{H} = \{h_{\lambda, \Phi} \mid \lambda \in \mathbb{R}_{\geq 0}^m, \Phi : Y \rightarrow \mathbb{R}^m\}$

$h_{\lambda, \Phi} : \text{SL} \rightarrow \text{TL}$

$$s \mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{y \in Y: \pi_{\text{SL}}(y)=s} \lambda \cdot \Phi(y) \right)$$

$$\lambda_1 \cdot \Phi(y)_1 + \dots + \lambda_m \cdot \Phi(y)_m$$

here:

$$\begin{aligned} m &= 2 \\ \Phi(y) &= (h_{\text{TM}}(y), h_{\text{LM}}(y)) \end{aligned}$$

h_{TM} : translation model

h_{LM} : language model

outline of the talk:

- ▶ Statistical machine translation
- ▶ Modeling with wta and wtt
- ▶ Using automata theoretic results to “improve” modeling
- ▶ Summary

first assumption:

SL and TL are the yields of
weighted recognizable tree languages

first assumption:

SL and TL are the yields of
weighted recognizable tree languages

weighted tree language: $L : T_\Sigma \rightarrow \mathbb{R}$

first assumption:

SL and TL are the yields of
weighted recognizable tree languages

weighted tree language: $L : T_\Sigma \rightarrow \mathbb{R}$

L is recognizable:

if there is a wta \mathcal{A}
which “recognizes” (computes) L

weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

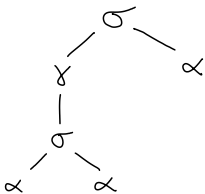
- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma)$ $\delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma)$ $\delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$

set of runs on ξ : $R_{\mathcal{A}}(\xi)$

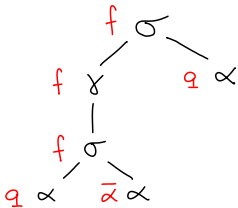


weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma) \quad \delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$

set of runs on ξ : $R_{\mathcal{A}}(\xi)$

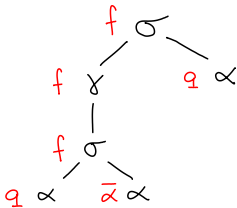


weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma) \quad \delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$ set of runs on ξ : $R_{\mathcal{A}}(\xi)$

weight of r : $\text{wt}(r) = \prod_{w \in \text{pos}(\xi)} \delta_\sigma(r(w_1) \cdots r(w_k), r(w))$



σ : label of ξ at w
 k : rank of σ

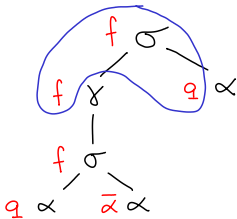
weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma) \quad \delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$ set of runs on ξ : $R_{\mathcal{A}}(\xi)$

weight of r : $\text{wt}(r) = \prod_{w \in \text{pos}(\xi)} \delta_\sigma(r(w_1) \cdots r(w_k), r(w))$

$w = \varepsilon$
 $\delta_\sigma(fq, f)$



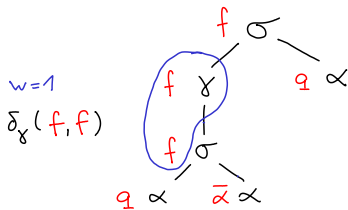
σ : label of ξ at w
 k : rank of σ

weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma) \quad \delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$ set of runs on ξ : $R_{\mathcal{A}}(\xi)$

weight of r : $\text{wt}(r) = \prod_{w \in \text{pos}(\xi)} \delta_\sigma(r(w_1) \cdots r(w_k), r(w))$

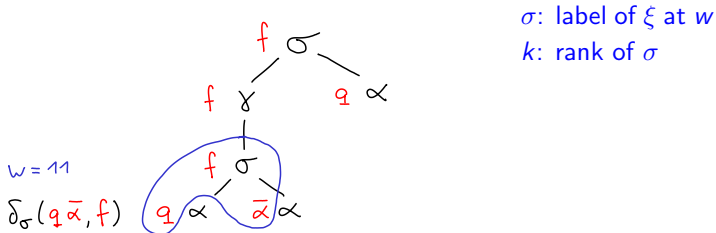


weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma) \quad \delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$ set of runs on ξ : $R_{\mathcal{A}}(\xi)$

weight of r : $\text{wt}(r) = \prod_{w \in \text{pos}(\xi)} \delta_\sigma(r(w_1) \cdots r(w_k), r(w))$

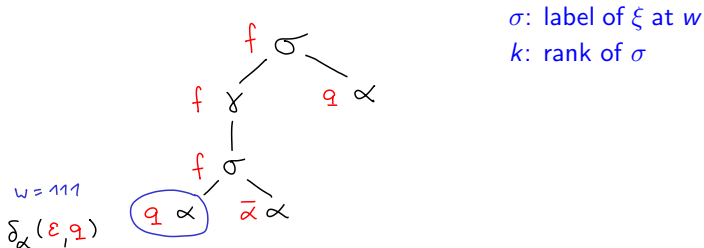


weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma) \quad \delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$ set of runs on ξ : $R_{\mathcal{A}}(\xi)$

weight of r : $\text{wt}(r) = \prod_{w \in \text{pos}(\xi)} \delta_\sigma(r(w_1) \cdots r(w_k), r(w))$



weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma)$ $\delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$ set of runs on ξ : $R_{\mathcal{A}}(\xi)$

weight of r : $\text{wt}(r) = \prod_{w \in \text{pos}(\xi)} \delta_\sigma(r(w1) \cdots r(wk), r(w))$

σ : label of ξ at w

k : rank of σ

weighted tree automaton (wta) $\mathcal{A} = (Q, \Sigma, \delta, F)$

- ▶ Q finite set (states)
- ▶ Σ ranked alphabet (input symbols)
- ▶ $\delta = (\delta_\sigma \mid \sigma \in \Sigma) \quad \delta_\sigma : Q^k \times Q \rightarrow \mathbb{R}$
 $\delta_\sigma(q_1 \cdots q_k, q) \in \mathbb{R}$
- ▶ $F \subseteq Q$ (final states)

run on $\xi \in T_\Sigma$: $r : \text{pos}(\xi) \rightarrow Q$ set of runs on ξ : $R_{\mathcal{A}}(\xi)$

weight of r : $\text{wt}(r) = \prod_{w \in \text{pos}(\xi)} \delta_\sigma(r(w1) \cdots r(wk), r(w))$

σ : label of ξ at w

k : rank of σ

weighted tree language recognized by \mathcal{A} :

$$L_{\mathcal{A}} : T_\Sigma \rightarrow \mathbb{R}, \quad L_{\mathcal{A}}(\xi) = \max_{r \in R_{\mathcal{A}}(\xi)} \text{wt}(r)$$

second assumption:

translation from SL and TL is specified by
a weighted tree transducer

[Yamada, Knight 01] translation from English to Japanese

weighted tree transducer (wtt) $\mathcal{M} = (Q, \Sigma, q_0, R)$

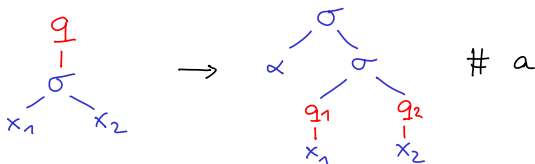
▶ Q, Σ as for wta

Σ : input and output symbols

▶ $q_0 \in Q$ (initial state)

weighted tree transducer (wtt) $\mathcal{M} = (Q, \Sigma, q_0, R)$

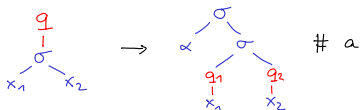
- ▶ Q, Σ as for wta Σ : input and output symbols
- ▶ $q_0 \in Q$ (initial state)
- ▶ R finite set of particular term rewrite rules with weights



linear, nondeleting in x_1, \dots, x_k

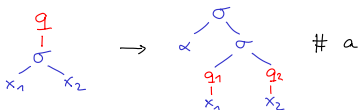
weighted tree transducer (wtt) $\mathcal{M} = (Q, \Sigma, q_0, R)$

- ▶ Q, Σ as for wta Σ : input and output symbols
- ▶ $q_0 \in Q$ (initial state)
- ▶ R finite set of particular term rewrite rules with weights



weighted tree transducer (wtt) $\mathcal{M} = (Q, \Sigma, q_0, R)$

- ▶ Q, Σ as for wta Σ : input and output symbols
- ▶ $q_0 \in Q$ (initial state)
- ▶ R finite set of particular term rewrite rules with weights

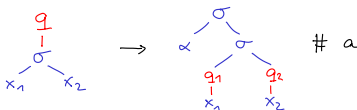


(leftmost) derivation:

$$d = \rho_1 \cdots \rho_n$$

weighted tree transducer (wtt) $\mathcal{M} = (Q, \Sigma, q_0, R)$

- ▶ Q, Σ as for wta Σ : input and output symbols
- ▶ $q_0 \in Q$ (initial state)
- ▶ R finite set of particular term rewrite rules with weights



(leftmost) derivation:

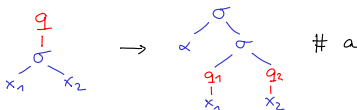
$$d = \rho_1 \cdots \rho_n$$

weight of a derivation d :

$$\text{wt}(d) = \prod_{i=1}^n \text{wt}(\rho_i)$$

weighted tree transducer (wtt) $\mathcal{M} = (Q, \Sigma, q_0, R)$

- ▶ Q, Σ as for wta Σ : input and output symbols
- ▶ $q_0 \in Q$ (initial state)
- ▶ R finite set of particular term rewrite rules with weights



(leftmost) derivation:

$$d = \rho_1 \cdots \rho_n$$

weight of a derivation d :

$$\text{wt}(d) = \prod_{i=1}^n \text{wt}(\rho_i)$$

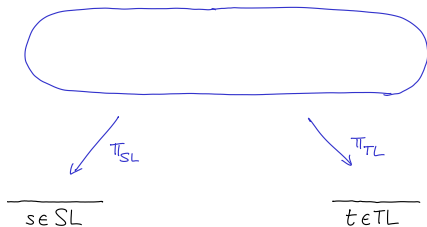
weighted tree transformation computed by \mathcal{M} :

$$\tau_{\mathcal{M}} : T_{\Sigma} \times T_{\Sigma} \rightarrow \mathbb{R}, \quad \tau_{\mathcal{M}}(\xi_1, \xi_2) = \max_{\substack{d \in D_{\mathcal{M}}: \\ \pi(d) = (\xi_1, \xi_2)}} \text{wt}(d)$$

$D_{\mathcal{M}}$: set of all derivations

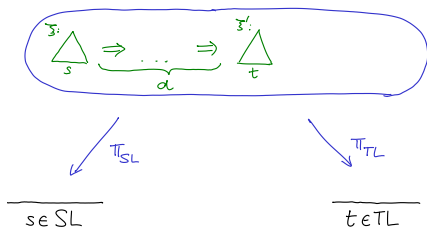
log-linear modeling with wtt and wta:

set Y of correspondence
structures [Liang et al. 06]



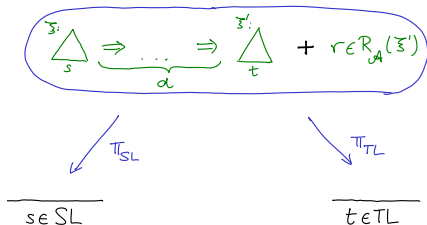
log-linear modeling with wtt and wta:

set Y of correspondence
structures [Liang et al. 06]



wtt \mathcal{M} as translation model; $d \in D_{\mathcal{M}}$

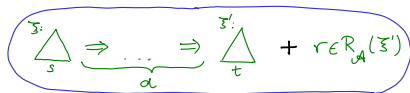
log-linear modeling with wtt and wta:



set Y of correspondence structures [Liang et al. 06]

wtt \mathcal{M} as translation model; $d \in D_{\mathcal{M}}$
wta \mathcal{A} as language model; $r \in R_{\mathcal{A}}$

log-linear modeling with wtt and wta:



set Y of correspondence structures [Liang et al. 06]



$s \in \text{SL}$



$t \in \text{TL}$

wtt \mathcal{M} as translation model; $d \in D_{\mathcal{M}}$
wta \mathcal{A} as language model; $r \in R_{\mathcal{A}}$

$$Y = \{(d, r) \in D_{\mathcal{M}} \times R_{\mathcal{A}} \mid r \in R_{\mathcal{A}}(\text{last}(d))\}$$

$$\pi_{\text{SL}}(d, r) = \text{yield}(\text{first}(d))$$

$$\pi_{\text{TL}}(d, r) = \text{yield}(\text{last}(d))$$

wtt \mathcal{M} as translation model; $d \in D_{\mathcal{M}}$

wta \mathcal{A} as language model; $r \in R_{\mathcal{A}}$

$$Y = \{(d, r) \in D_{\mathcal{M}} \times R_{\mathcal{A}} \mid r \in R_{\mathcal{A}}(\text{last}(d))\}$$

wtt \mathcal{M} as translation model; $d \in D_{\mathcal{M}}$

wta \mathcal{A} as language model; $r \in R_{\mathcal{A}}$

$$Y = \{(d, r) \in D_{\mathcal{M}} \times R_{\mathcal{A}} \mid r \in R_{\mathcal{A}}(\text{last}(d))\}$$

here:

$$m = 2$$
$$\Phi(d, r) = \left(\log \text{wt}_{\mathcal{M}}(d), \log \text{wt}_{\mathcal{A}}(r) \right)$$

wtt \mathcal{M} as translation model; $d \in D_{\mathcal{M}}$

wta \mathcal{A} as language model; $r \in R_{\mathcal{A}}$

$$Y = \{(d, r) \in D_{\mathcal{M}} \times R_{\mathcal{A}} \mid r \in R_{\mathcal{A}}(\text{last}(d))\}$$

here:

$$m = 2$$
$$\Phi(d, r) = \left(\log \text{wt}_{\mathcal{M}}(d), \log \text{wt}_{\mathcal{A}}(r) \right)$$

$$\mathcal{H} = \{h_{\lambda, \mathcal{M}, \mathcal{A}} \mid \lambda \in \mathbb{R}_{\geq 0}^2, \text{ wtt } \mathcal{M}, \text{ wta } \mathcal{A}\}$$

$$h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} \rightarrow \text{TL}$$

$$s \mapsto \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d, r) = s}{\text{argmax}}_{(d, r) \in Y} \text{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2} \right)$$

outline of the talk:

- ▶ Statistical machine translation
- ▶ Modeling with wta and wtt
- ▶ Using automata theoretic results to “improve” modeling
 - ▶ Weight exponentiation
 - ▶ Output product
- ▶ Summary

$$h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} \rightarrow \text{TL}$$
$$s \mapsto \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d,r)=s}{\operatorname{argmax}}_{(d,r) \in Y} \operatorname{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \operatorname{wt}_{\mathcal{A}}(r)^{\lambda_2} \right)$$

$$\begin{aligned}
h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} &\rightarrow \text{TL} \\
s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{\pi_{\text{SL}}(d,r)=s} (d,r) \in \mathcal{Y} : \text{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2} \right) \\
&= \pi_{\text{TL}} \left(\operatorname{argmax}_{\pi_{\text{SL}}(d,r)=s} (d,r) \in \mathcal{Y} : \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2/\lambda_1} \right)
\end{aligned}$$

$$\begin{aligned}
h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} &\rightarrow \text{TL} \\
s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{\pi_{\text{SL}}(d,r)=s} (d,r) \in \mathcal{Y} : \text{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2} \right) \\
&= \pi_{\text{TL}} \left(\operatorname{argmax}_{\pi_{\text{SL}}(d,r)=s} (d,r) \in \mathcal{Y} : \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2/\lambda_1} \right)
\end{aligned}$$

Lemma: Let \mathcal{A} be a wta and $\lambda \in \mathbb{R}_{\geq 0}$.

$$\begin{aligned}
h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} &\rightarrow \text{TL} \\
s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{\pi_{\text{SL}}(d,r)=s} (d,r) \in \mathcal{Y} : \text{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2} \right) \\
&= \pi_{\text{TL}} \left(\operatorname{argmax}_{\pi_{\text{SL}}(d,r)=s} (d,r) \in \mathcal{Y} : \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2/\lambda_1} \right)
\end{aligned}$$

Lemma: Let \mathcal{A} be a wta and $\lambda \in \mathbb{R}_{\geq 0}$.

There is a wta \mathcal{A}' s.t. $Q_{\mathcal{A}'} = Q_{\mathcal{A}}$ and
 $\text{wt}_{\mathcal{A}'}(r) = \text{wt}_{\mathcal{A}}(r)^{\lambda}$ for every r .

$$\begin{aligned}
h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} &\rightarrow \text{TL} \\
s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{(d,r) \in \mathcal{Y} \\ \pi_{\text{SL}}(d,r)=s}} \text{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2} \right) \\
&= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{(d,r) \in \mathcal{Y} \\ \pi_{\text{SL}}(d,r)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2/\lambda_1} \right)
\end{aligned}$$

Lemma: Let \mathcal{A} be a wta and $\lambda \in \mathbb{R}_{\geq 0}$.

There is a wta \mathcal{A}' s.t. $Q_{\mathcal{A}'} = Q_{\mathcal{A}}$ and
 $\text{wt}_{\mathcal{A}'}(r) = \text{wt}_{\mathcal{A}}(r)^\lambda$ for every r .

$$= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{(d,r) \in \mathcal{Y}' \\ \pi_{\text{SL}}(d,r)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}'}(r) \right)$$

outline of the talk:

- ▶ Statistical machine translation
- ▶ Modeling with wta and wtt
- ▶ Using automata theoretic results to “improve” modeling
 - ▶ Weight exponentiation
 - ▶ Output product
- ▶ Summary

Let $\tau : T_\Sigma \times T_\Sigma \rightarrow \mathbb{R}$ and $L : T_\Sigma \rightarrow \mathbb{R}$

Let $\tau : T_\Sigma \times T_\Sigma \rightarrow \mathbb{R}$ and $L : T_\Sigma \rightarrow \mathbb{R}$

output product of τ and L :

$$\tau \triangleright L : T_\Sigma \times T_\Sigma \rightarrow \mathbb{R}$$

$$(\tau \triangleright L)(\xi, \zeta) \mapsto \tau(\xi, \zeta) \cdot L(\zeta)$$

Let $\tau : T_\Sigma \times T_\Sigma \rightarrow \mathbb{R}$ and $L : T_\Sigma \rightarrow \mathbb{R}$

output product of τ and L :

$$\tau \triangleright L : T_\Sigma \times T_\Sigma \rightarrow \mathbb{R}$$

$$(\tau \triangleright L)(\xi, \zeta) \mapsto \tau(\xi, \zeta) \cdot L(\zeta)$$

input product of L and τ :

$$L \triangleleft \tau : T_\Sigma \times T_\Sigma \rightarrow \mathbb{R}$$

$$(L \triangleleft \tau)(\xi, \zeta) \mapsto L(\xi) \cdot \tau(\xi, \zeta)$$

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

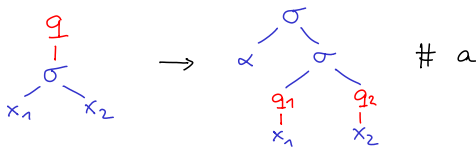
Proof: [Baker 79, Engelfriet, Fülöp, V. 02]

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

Proof: [Baker 79, Engelfriet, Fülöp, V. 02]

rule of \mathcal{M} :

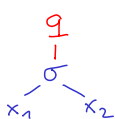


Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

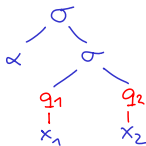
There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

Proof: [Baker 79, Engelfriet, Fülöp, V. 02]

rule of \mathcal{M} :



→



a

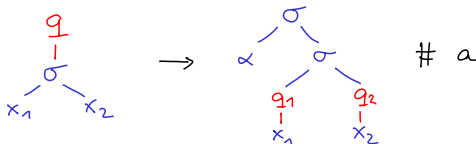
states of \mathcal{A} : p, p_1, p_2

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

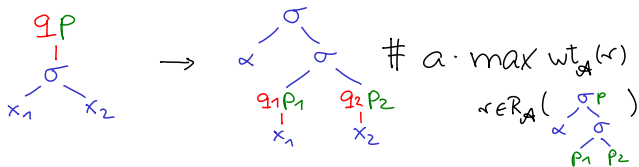
Proof: [Baker 79, Engelfriet, Fülöp, V. 02]

rule of \mathcal{M} :



states of \mathcal{A} : p, p_1, p_2

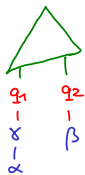
rule of $\mathcal{M} \triangleright \mathcal{A}$:



d:



\Rightarrow



\Rightarrow



\Rightarrow



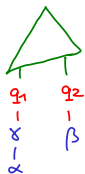
\Rightarrow



d:



\Rightarrow



\Rightarrow



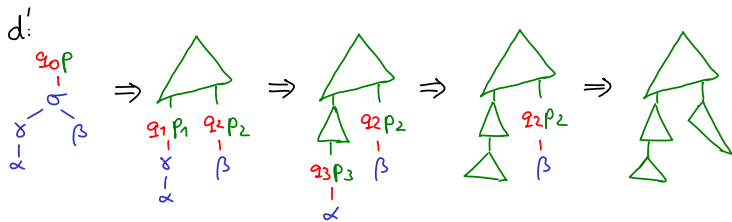
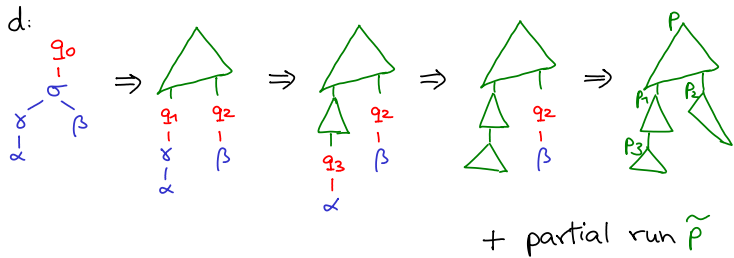
\Rightarrow

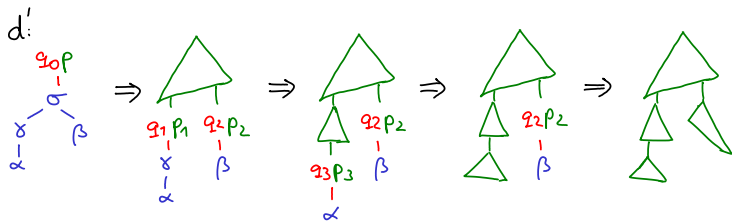
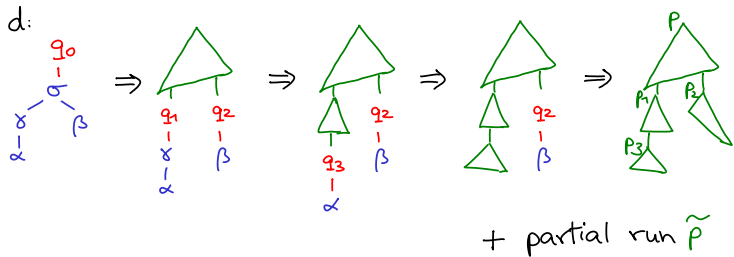


\Rightarrow



+ partial run \tilde{P}





$\varphi : D_{\mathcal{M}'} \rightarrow \{(d, \tilde{p}) \mid d \in D_{\mathcal{M}}, \tilde{p} \in R_A^{\text{partial}}(d)\}$ bijection

$\text{wt}(d') = \text{wt}(d) \cdot \max_{r \in \text{completion}(\tilde{p})} \text{wt}(r)$

recall:

$$h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} \rightarrow \text{TL}$$

$$s \mapsto \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d,r)=s}{\operatorname{argmax}}_{(d,r) \in \mathcal{Y}'} \operatorname{wt}_{\mathcal{M}}(d) \cdot \operatorname{wt}_{\mathcal{A}'}(r) \right)$$

$h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} \rightarrow \text{TL}$

$$\begin{aligned} s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{(d,r) \in \mathcal{Y}' \\ \pi_{\text{SL}}(d,r)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}'}(r) \right) \\ &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M}} \\ \pi_{\text{SL}}(d)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \max_{r \in R_{\mathcal{A}'}(\text{last}(d))} \text{wt}_{\mathcal{A}'}(r) \right) \end{aligned}$$

$h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} \rightarrow \text{TL}$

$$\begin{aligned} s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{(d,r) \in \mathcal{Y}'} : \pi_{\text{SL}}(d,r)=s} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}'}(r) \right) \\ &= \pi_{\text{TL}} \left(\operatorname{argmax}_{d \in D_{\mathcal{M}}} : \pi_{\text{SL}}(d)=s} \text{wt}_{\mathcal{M}}(d) \cdot \max_{r \in R_{\mathcal{A}'}(\text{last}(d))} \text{wt}_{\mathcal{A}'}(r) \right) \\ &= \pi_{\text{TL}} \left(\operatorname{argmax}_{d \in D_{\mathcal{M}}} : \pi_{\text{SL}}(d)=s} \text{wt}_{\mathcal{M}}(d) \cdot \max_{\tilde{p} \in R_{\mathcal{A}'}^{\text{P}}(d)} \max_{r \in \text{compl.}(\tilde{p})} \text{wt}_{\mathcal{A}'}(r) \right) \end{aligned}$$

$$h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} \rightarrow \text{TL}$$

$$\begin{aligned}
 s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{(d,r) \in \mathcal{Y}' \\ \pi_{\text{SL}}(d,r)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}'}(r) \right) \\
 &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M}} \\ \pi_{\text{SL}}(d)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \max_{r \in R_{\mathcal{A}'}(\text{last}(d))} \text{wt}_{\mathcal{A}'}(r) \right) \\
 &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M}} \\ \pi_{\text{SL}}(d)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \max_{\tilde{p} \in R_{\mathcal{A}'}^{\text{P}}(d)} \max_{r \in \text{compl.}(\tilde{p})} \text{wt}_{\mathcal{A}'}(r) \right) \\
 &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M}} \\ \pi_{\text{SL}}(d)=s \\ \tilde{p} \in R_{\mathcal{A}'}^{\text{P}}(d)}} \text{wt}_{\mathcal{M}}(d) \cdot \max_{r \in \text{compl.}(\tilde{p})} \text{wt}_{\mathcal{A}'}(r) \right)
 \end{aligned}$$

$$h_{\lambda, \mathcal{M}, \mathcal{A}} : \text{SL} \rightarrow \text{TL}$$

$$\begin{aligned}
 s &\mapsto \pi_{\text{TL}} \left(\operatorname{argmax}_{(d,r) \in \mathcal{Y}'} : \pi_{\text{SL}}(d,r)=s} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}'}(r) \right) \\
 &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M}} \\ \pi_{\text{SL}}(d)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \max_{r \in R_{\mathcal{A}'}(\text{last}(d))} \text{wt}_{\mathcal{A}'}(r) \right) \\
 &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M}} \\ \pi_{\text{SL}}(d)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \max_{\tilde{p} \in R_{\mathcal{A}'}^{\text{P}}(d)} \max_{r \in \text{compl.}(\tilde{p})} \text{wt}_{\mathcal{A}'}(r) \right) \\
 &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M}} \\ \pi_{\text{SL}}(d)=s \\ \tilde{p} \in R_{\mathcal{A}'}^{\text{P}}(d)}} \text{wt}_{\mathcal{M}}(d) \cdot \max_{r \in \text{compl.}(\tilde{p})} \text{wt}_{\mathcal{A}'}(r) \right) \\
 &= \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M} \triangleright \mathcal{A}'} \\ \pi_{\text{SL}}(d)=s}} \text{wt}_{\mathcal{M} \triangleright \mathcal{A}'}(d) \right)
 \end{aligned}$$

recall:

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

recall:

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

generalization to mildly context-sensitive languages

Theorem [Büchse, Nederhof, V. 11]:

Let \mathcal{M} synchronized tree-adjoining grammar (STAG)
and \mathcal{A} wta.

There is an STAG $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

recall:

Theorem [Maletti 06]: Let \mathcal{M} wtt and \mathcal{A} wta.

There is a wtt $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

generalization to mildly context-sensitive languages

Theorem [Büchse, Nederhof, V. 11]:

Let \mathcal{M} synchronized tree-adjoining grammar (STAG)
and \mathcal{A} wta.

There is an STAG $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

Theorem [Nederhof, V. 12]:

Let \mathcal{M} synchronized context-free tree grammar (SCFTG)
and \mathcal{A} wta.

There is an SCFTG $\mathcal{M} \triangleright \mathcal{A}$ such that: $\tau_{\mathcal{M} \triangleright \mathcal{A}} = \tau_{\mathcal{M}} \triangleright L_{\mathcal{A}}$

outline of the talk:

- ▶ Statistical machine translation
- ▶ Modeling with wta and wtt
- ▶ Using automata theoretic results to “improve” modeling
- ▶ **Summary**

- ▶ model for translation from SL to TL: $w_{tt} \mathcal{M}$
- ▶ model for TL: $w_{ta} \mathcal{A}$
- ▶ sentence s of SL

- ▶ model for translation from SL to TL: wtt \mathcal{M}
- ▶ model for TL: wta \mathcal{A}
- ▶ sentence s of SL

$$h_{\lambda, \mathcal{M}, \mathcal{A}}(s) = \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d,r)=s}{\operatorname{argmax}}_{(d,r) \in \mathcal{Y}}: \operatorname{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \operatorname{wt}_{\mathcal{A}}(r)^{\lambda_2} \right)$$

- ▶ model for translation from SL to TL: wtt \mathcal{M}
- ▶ model for TL: wta \mathcal{A}
- ▶ sentence s of SL

$$h_{\lambda, \mathcal{M}, \mathcal{A}}(s) = \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d,r)=s}{\operatorname{argmax}}_{(d,r) \in \mathcal{Y}}: \operatorname{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \operatorname{wt}_{\mathcal{A}}(r)^{\lambda_2} \right)$$

$$\text{(weight exp.)} = \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d,r)=s}{\operatorname{argmax}}_{(d,r) \in \mathcal{Y}'}: \operatorname{wt}_{\mathcal{M}}(d) \cdot \operatorname{wt}_{\mathcal{A}'}(r) \right)$$

- ▶ model for translation from SL to TL: wtt \mathcal{M}
- ▶ model for TL: wta \mathcal{A}
- ▶ sentence s of SL

$$h_{\lambda, \mathcal{M}, \mathcal{A}}(s) = \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{(d,r) \in \mathcal{Y}: \\ \pi_{\text{SL}}(d,r)=s}} \text{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2} \right)$$

$$\text{(weight exp.)} = \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{(d,r) \in \mathcal{Y}': \\ \pi_{\text{SL}}(d,r)=s}} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}'}(r) \right)$$

$$\text{(output product)} = \pi_{\text{TL}} \left(\operatorname{argmax}_{\substack{d \in D_{\mathcal{M} \triangleright \mathcal{A}'}: \\ \pi_{\text{SL}}(d)=s}} \text{wt}_{\mathcal{M} \triangleright \mathcal{A}'}(d) \right)$$

- ▶ model for translation from SL to TL: wtt \mathcal{M}
- ▶ model for TL: wta \mathcal{A}
- ▶ sentence s of SL

$$h_{\lambda, \mathcal{M}, \mathcal{A}}(s) = \pi_{\text{TL}} \left(\operatorname{argmax}_{(d,r) \in Y: \pi_{\text{SL}}(d,r)=s} \text{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \text{wt}_{\mathcal{A}}(r)^{\lambda_2} \right)$$

$$\text{(weight exp.)} = \pi_{\text{TL}} \left(\operatorname{argmax}_{(d,r) \in Y': \pi_{\text{SL}}(d,r)=s} \text{wt}_{\mathcal{M}}(d) \cdot \text{wt}_{\mathcal{A}'}(r) \right)$$

$$\text{(output product)} = \pi_{\text{TL}} \left(\operatorname{argmax}_{d \in D_{\mathcal{M} \triangleright \mathcal{A}'}: \pi_{\text{SL}}(d)=s} \text{wt}_{\mathcal{M} \triangleright \mathcal{A}'}(d) \right)$$

$$\text{(B,S,P; input product)} = \pi_{\text{TL}} \left(\operatorname{argmax}_{d \in D_{\mathcal{A}_s \triangleleft (\mathcal{M} \triangleright \mathcal{A})} \text{wt}_{\mathcal{A}_s \triangleleft (\mathcal{M} \triangleright \mathcal{A}')} (d) \right)$$

- ▶ model for translation from SL to TL: wtt \mathcal{M}
- ▶ model for TL: wta \mathcal{A}
- ▶ sentence s of SL

$$h_{\lambda, \mathcal{M}, \mathcal{A}}(s) = \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d,r)=s}{\operatorname{argmax}}_{(d,r) \in Y'}: \operatorname{wt}_{\mathcal{M}}(d)^{\lambda_1} \cdot \operatorname{wt}_{\mathcal{A}}(r)^{\lambda_2} \right)$$

$$\text{(weight exp.)} = \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d,r)=s}{\operatorname{argmax}}_{(d,r) \in Y'}: \operatorname{wt}_{\mathcal{M}}(d) \cdot \operatorname{wt}_{\mathcal{A}'}(r) \right)$$

$$\text{(output product)} = \pi_{\text{TL}} \left(\underset{\pi_{\text{SL}}(d)=s}{\operatorname{argmax}}_{d \in D_{\mathcal{M} \triangleright \mathcal{A}'}}: \operatorname{wt}_{\mathcal{M} \triangleright \mathcal{A}'}(d) \right)$$

$$\begin{aligned} \text{(B,S,P; input product)} &= \pi_{\text{TL}} \left(\underset{d \in D_{\mathcal{A}_s \triangleleft (\mathcal{M} \triangleright \mathcal{A}')}}{\operatorname{argmax}}: \operatorname{wt}_{\mathcal{A}_s \triangleleft (\mathcal{M} \triangleright \mathcal{A}')} (d) \right) \\ &= \pi_{\text{TL}} \left(\operatorname{Knuth}(\mathcal{A}_s \triangleleft (\mathcal{M} \triangleright \mathcal{A}')) \right) \end{aligned}$$

impression: SMT is “playing with formulas”

impression: SMT is “playing with formulas”

but: SMT is an engineering task!

impression: SMT is “playing with formulas”

but: SMT is an engineering task!

weighted tree automata and weighted tree transducers

can help in **modeling**

statistical machine translation of natural languages

References:

- ▶ [Baker 79] Composition of top-down and bottom-up tree transductions
- ▶ [Bar-Hillel, Shamir, Perles 61] On formal properties of simple phrase structure grammars
- ▶ [Büchse, Nederhof, V. 11] Tree Parsing with Synchronous Tree-Adjoining Grammars
- ▶ [Engelfriet, Fülöp, V. 02] Bottom-up and Top-down Tree Series Transformations
- ▶ [Fülöp, V. 09] Weighted tree automata and tree transducers
- ▶ [Knight et al. 03-...] ...
- ▶ [Knuth 77] A generalization of Dijkstra's algorithm
- ▶ [Lopez 08] Statistical Machine Translation
- ▶ [Liang, Bouchard-Côté, Klein, Taskar 06] An End-to-End Discriminative Approach to Machine Translation
- ▶ [Maletti 06] Compositions of Tree Series Transformations
- ▶ [Maletti, Satta 09] Parsing Algorithms based on Tree Automata
- ▶ [Nederhof, V. 12] Synchronous Context-Free Tree Grammars
- ▶ [Thatcher 67] Characterizing derivation trees of context-free grammars through a generalization of finite automata theory
- ▶ [Yamada, Knight 01] A syntax-based statistical translation model

[Knight et al. 03-...]

- ▶ [Charniak, Knight, Yamada 03] Syntax-based language models for statistical machine translation
- ▶ [Galley, Hopkins, Knight, Marcu 04] What's in a translation rule?
- ▶ [Graehl, Knight 04] Training tree transducers
- ▶ [Graehl, Knight 05] An overview of probabilistic tree transducers for natural language processing

[Knight et al. 03-...]

- ▶ [Charniak, Knight, Yamada 03] Syntax-based language models for statistical machine translation
- ▶ [Galley, Hopkins, Knight, Marcu 04] What's in a translation rule?
- ▶ [Graehl, Knight 04] Training tree transducers
- ▶ [Graehl, Knight 05] An overview of probabilistic tree transducers for natural language processing

[Maletti 11]

Survey: Weighted Extended Top-down Tree Transducers –
Part III: Applications in Machine Translation

[Knight et al. 03-...]

- ▶ [Charniak, Knight, Yamada 03] Syntax-based language models for statistical machine translation
- ▶ [Galley, Hopkins, Knight, Marcu 04] What's in a translation rule?
- ▶ [Graehl, Knight 04] Training tree transducers
- ▶ [Graehl, Knight 05] An overview of probabilistic tree transducers for natural language processing

[Maletti 11]

Survey: Weighted Extended Top-down Tree Transducers –
Part III: Applications in Machine Translation

Thank you!