### A Shared Task on Bandit Learning for Machine Translation

<u>Artem Sokolov</u><sup>◊</sup>, Julia Kreutzer\*, Kellen Sunderland<sup>◊</sup>, Pavel Danchenko<sup>◊</sup>, Witold Szymaniak<sup>◊</sup>, Hagen Fürstenau<sup>◊</sup>, Stefan Riezler<sup>\*,‡</sup>



WMT17

MT supervision	cost	noise/signal	
references	<b>\$\$\$</b>	low	
post-edits/quasi-refs	\$\$	middle	
feedback	\$	high	

# Motivation

MT supervision	cost	noise/signal	downstream fit
references	<b>\$\$\$</b>	low	hard
post-edits/quasi-refs	\$\$	middle	hard
feedback	\$	high	easy/trivial

## Motivation

MT supervision	cost	noise/signal	downstream fit
references	\$\$\$	low	hard
post-edits/quasi-refs	\$\$	middle	hard
feedback	\$	high	easy/trivial

#### Sweet spot for bandit MT:

- **1** costs drop faster than noise/signal increases  $\Rightarrow$  lots of data
- 2 downstream performance is hard to wrap into an automatic metric

One-armed bandits - slot-machines:

- pull an arm to play
- get some reward (or none)
- try a new machine or stick to the best among discovered so far?



#### One-armed Bandit

Multi-armed bandits:

- many arms (actions)
- each arm has an unknown reward distribution
- find an arm-picking strategy to maximize total reward



Multi-armed Bandit

Multi-armed bandits for structured prediction (MT):

- observe context (source sentence)
- pick one out of exponentially many outputs (translations)
- each structure results in some reward (BLEU)
- tune an arm-picking strategy (decoder weights)



Exponential number of arms

Multi-armed bandits for structured prediction (MT):

- observe context (source sentence)
- pick one out of exponentially many outputs (translations)
- each structure results in some reward (BLEU)
- tune an arm-picking strategy (decoder weights)



MT - Full info

Multi-armed bandits for structured prediction (MT):

- observe context (source sentence)
- pick one out of exponentially many outputs (translations)
- each structure results in some reward (BLEU)
- tune an arm-picking strategy (decoder weights)
- note: only one translation is scored, others are not



MT - Partial info

#### Shared Task



for  $t = 0, \ldots, T$  do

Request source sentence  $x_t$  from service Propose a translation  $y_t$ Obtain feedback  $\Delta(y_t)$  from service Improve MT model

- DE-EN (pre-processed)
- domain-adaptation: general (WMT17)  $\rightarrow$  e-commerce (Amazon)
- all participants received the same sequence of  $x_t$
- feedback was sent-BLEU (plans for more realistic feedback dropped to avoid complicating the task)
- organizers provided the service, client SDK and baselines

- DE-EN (pre-processed)
- domain-adaptation: general (WMT17) → e-commerce (Amazon)
- all participants received the same sequence of  $x_t$
- feedback was sent-BLEU (plans for more realistic feedback dropped to avoid complicating the task)
- organizers provided the service, client SDK and baselines

#### Participants had to do:

- 1 pick Python or Java
- 2 download a short client snippet
- **3** wrap it around an MT system

- DE-EN (pre-processed)
- domain-adaptation: general (WMT17) → e-commerce (Amazon)
- all participants received the same sequence of  $x_t$
- feedback was sent-BLEU (plans for more realistic feedback dropped to avoid complicating the task)
- organizers provided the service, client SDK and baselines

#### Participants had to do:

- 1 pick Python or Java
- 2 download a short client snippet
- 3 wrap it around an MT system

phase	sentences	passes	purpose
mock (since 13 Mar)	40	unlimited	test client API
dev (since 5 Apr)	40k	unlimited	tune hyperparams
train (25 Apr - Jun 9)	1.3M	only one	final evaluation

cumulative reward

$$\sum_{t=1}^{T} \Delta(y_t)$$

**corpus-BLEU** at regular intervals on an embedded test set:

- ➡ 700 sent. at 4 locations in 40k dev sent.
- ➡ 4000 sent. at 12 locations in 1.3M train sent.

regret

$$\frac{1}{T}\sum_{t=1}^{T}\Delta(y_t^*) - \Delta(y_t)$$

(average cumulative reward difference w.r.t. to an in-domain system)

- different domains ⇒ high OOV rate
- learning from one-shot feedback
- real-world data:
  - ➡ typos/errors in sources
  - mixed direction of data
  - translators improved readability/corrected errors/deleted irrelevancies

- different domains  $\Rightarrow$  high OOV rate
- learning from one-shot feedback
- real-world data:
  - ➡ typos/errors in sources
  - mixed direction of data
  - translators improved readability/corrected errors/deleted irrelevancies

source	reference
schwarz gr.xxl / xxxl , 147 cm	black <mark>, size</mark> xxl / xxxl 147 cm
für starke , glänzende nägel	great for strengthen your nails and enhance shine
seemless verarbeitung maschinenwaschbar bei 30 ° c	seamless processing machine washable at 30 degrees .
32_unzen volumen material : 1050 denier nylon . für e-gitarre entworfen	32-ounce capacity material : 1050d nylon . designed for electric guitar

- 8 teams registered
- 4 used the dev service
- 2 started full training:



Guillaume Wisniewski

- ➡ UCB-style selection from a pool of MT systems
- online linear regression to predict rewards
- additional exploration



Amr Sharaf, Shi Feng, Khanh Nguyen, Kianté Brantley, Hal Daumé

- domain-adaptation (Moore-Lewis)
- online non-linear regression for rewards
- policy gradient with adaptive control variate (Advantage Actor-Critic)

gradient-free

- NMT (word- and BPE-based; NeuralMonkey)
- SMT (dense features; cdec)

stochastic approximation of the expected loss

### policy-gradient

$$\begin{split} \mathbb{E}_{\tilde{y} \sim p_w(y|x)}[\Delta(\tilde{y})] \\ \text{for } t = 1, \dots, T \text{ do} \\ \text{Observe } x_t \\ \text{Sample } \tilde{y}_t \sim p_{w_t}(y|x_t) \\ \text{Obtain } \Delta(\tilde{y}_t) \end{split}$$

$$w_{t+1} = w_t - \gamma \Delta(\tilde{y}_t) \nabla \log p_{w_t}(\tilde{y}_t | x_t)$$

 $\mathbb{E}_{\varepsilon \sim \mathcal{N}(0,1)}[\Delta(\hat{y}(w + \varepsilon))]$  **for**  $t = 1, \dots, T$  **do** Observe  $x_t$ Sample  $\varepsilon_t \sim \mathcal{N}(0, 1)$ Decode  $\hat{y}_t$  with  $w_t + \varepsilon_t$ Obtain  $\Delta(\hat{y}_t)$  $w_{t+1} = w_t + \gamma \Delta(\hat{y}_t)\varepsilon_t$ 













# Checkpoints (corpus-BLEU)





#### Static system:

★ UMD domain adaptation got the best BLEU and lowest regret

#### Static system:

★ UMD domain adaptation got the best BLEU and lowest regret

#### Learning systems:

- cumulative reward:
  - ★ BNMT-EL is the only to beat its static NMT baseline
    - +lowest regret of all learning systems
    - +best sent-BLEU overall
  - ★ SMT-EL-CV came very close
- corpus-BLEU:
  - ★ SMT-EL-CV improves over its SMT baseline
  - none of the submissions show monotonic learning curves (or are too short)

## Static system:

★ UMD domain adaptation got the best BLEU and lowest regret

#### Learning systems:

- cumulative reward:
  - ★ BNMT-EL is the only to beat its static NMT baseline
    - +lowest regret of all learning systems
    - +best sent-BLEU overall
  - ★ SMT-EL-CV came very close
- corpus-BLEU:
  - ★ SMT-EL-CV improves over its SMT baseline
  - none of the submissions show monotonic learning curves (or are too short)

#### **Conclusion:**

- difficult task (even with all simplifications)
- difficult data
- we need more research!:)





task idea, NMT baselines





API, SDK, leaderboard & operation

task design, data, SMT baselines, etc.





task idea, NMT baselines and bandit pictures!

API, SDK, leaderboard & operation







task design, data, SMT baselines, etc.