Imitation Attacks and Defenses for Black-box Machine Translation Systems

Eric Wallace, Mitchell Stern, Dawn Song





Eric Wallace



Mitchell Stern



Dawn Song

Production NLP Models Are Lucrative

Production NLP Models Are Lucrative





Machine Translation



Smart Assistants

Production NLP Models Are Lucrative





Machine Translation



Smart Assistants

Result of large investments into data annotation and model design



Fake News Detection



Dialogue Systems



Machine Translation



Spam Filtering



Fake News Detection



Dialogue Systems



Machine Translation



Spam Filtering

Errors can have **negative societal consequences**



Errors can have **negative societal consequences**

An adversary can benefit financially by **stealing models**

An adversary can benefit financially by **stealing models**

• avoid long-term API costs by stealing models upfront

An adversary can benefit financially by **stealing models**

- avoid long-term API costs by stealing models upfront
- launch a competitor service of similar quality

An adversary can benefit financially by **stealing models**

- avoid long-term API costs by stealing models upfront
- launch a competitor service of similar quality

An adversary can benefit financially or harm society by **breaking models**

An adversary can benefit financially by **stealing models**

- avoid long-term API costs by stealing models upfront
- launch a competitor service of similar quality

An adversary can benefit financially or harm society by **breaking models**

• manipulate the stock market by fooling sentiment models

An adversary can benefit financially by **stealing models**

- avoid long-term API costs by stealing models upfront
- launch a competitor service of similar quality

An adversary can benefit financially or harm society by **breaking models**

- manipulate the stock market by fooling sentiment models
- bypass classifiers of fake news or hate speech

• Common Practice: keep data + model hidden

• Common Practice: keep data + model hidden



Hidden Data + Model

• Common Practice: keep data + model hidden



- Common Practice: keep data + model hidden
- Our paper: this is **not enough** to protect NLP models!



- Common Practice: keep data + model hidden
- Our paper: this is **not enough** to protect NLP models!
 - o adversaries can imitate black-box models



- Common Practice: keep data + model hidden
- Our paper: this is **not enough** to protect NLP models!
 - adversaries can imitate black-box models
 - imitation models help break black-box models



- Common Practice: keep data + model hidden
- Our paper: this is **not enough** to protect NLP models!
 - adversaries can imitate black-box models
 - imitation models help break black-box models
 - new defenses mitigate adversaries



- Common Practice: keep data + model hidden
- Our paper: this is not enough to protect NLP models!
 - adversaries can imitate black-box models
 - imitation models help break black-box models
 - new defenses mitigate adversaries
- We consider machine translation (MT) as a case study



• Goal: train imitation model that is similar to black-box API

- Goal: train imitation model that is similar to black-box API
- Method: query sentences and use API output as training data

- Goal: train imitation model that is similar to black-box API
- Method: query sentences and use API output as training data
- Not just model distillation:
 - \circ unknown data distribution

- Goal: train imitation model that is similar to black-box API
- Method: query sentences and use API output as training data
- Not just model distillation:
 - unknown data distribution
 - no distribution or feature matching losses

Setup:

• Black-box MT victim model for German-English

Setup:

- Black-box MT victim model for German-English
- Vary imitation model's architecture and queried sentences

Setup:

- Black-box MT victim model for German-English
- Vary imitation model's architecture and queried sentences

Evaluation metrics:

• BLEU on in-domain and out-of-domain data

Setup:

- Black-box MT victim model for German-English
- Vary imitation model's architecture and queried sentences

Evaluation metrics:

- BLEU on in-domain and out-of-domain data
- Output similarity using inter-system BLEU

Setup:

- Black-box MT victim model for German-English
- Vary imitation model's architecture and queried sentences

Evaluation metrics:

- BLEU on in-domain and out-of-domain data
- Output similarity using inter-system BLEU

For all architectures, data settings, and evaluation metrics, the imitation models closely match their victims

Imitating Production Models

• Imitate production systems on English-German and Nepali-English
Imitating Production Models

• Imitate production systems on English-German and Nepali-English

• We closely match production systems

	Model	Google	Bing	Systran
In-domain BLEU	Official	32.0	32.9	27.8
	Imitation	31.5	32.4	27.6

Imitating Production Models

• Imitate production systems on English-German and Nepali-English

• We closely match production systems

	Model	Google	Bing	Systran
In-domain BLEU	Official	32.0	32.9	27.8
	Imitation	31.5	32.4	27.6
Out-of-domain BLEU	Official	32.0	32.7	32.0
	Imitation	31.1	32.0	31.4

Most adversarial attacks for NLP assume white-box access
 How to do black-box attacks?

- Most adversarial attacks for NLP assume white-box access
 o How to do black-box attacks?
- Simple idea: transfer attacks from imitation models

- Most adversarial attacks for NLP assume white-box access
 o How to do black-box attacks?
- Simple idea: transfer attacks from imitation models

es ist über 7°F

German Decoder English Encoder it's over 7° F

- Most adversarial attacks for NLP assume white-box access
 o How to do black-box attacks?
- Simple idea: transfer attacks from imitation models



- Most adversarial attacks for NLP assume white-box access
 o How to do black-box attacks?
- Simple idea: transfer attacks from imitation models



- Most adversarial attacks for NLP assume white-box access
 o How to do black-box attacks?
- Simple idea: transfer attacks from imitation models



- Most adversarial attacks for NLP assume white-box access
 o How to do black-box attacks?
- Simple idea: transfer attacks from imitation models



We consider benign but representative adversarial attacks
 o same techniques would apply to malicious use cases

Untargeted Universal Trigger Attack

• A phrase that commonly causes mistranslations when appended to any input

Untargeted Universal Trigger Attack

• A phrase that commonly causes mistranslations when appended to any input

SYSTRAN

l heard machine translation is now superhuman Siehe Siehe Siehe Siehe Siehe Siehe Siehe In diesem Jahr ist es wieder soweit: Manche Manuskripte haben sich in der Hauptsache in der Hauptsache wieder in den Vordergrund gestellt.

(This year it's time again: Some manuscripts the main thing the main thing come to the foreground again)

Untargeted Universal Trigger Attack

• A phrase that commonly causes mistranslations when appended to any input

SYSTRAN

I heard machine translation is now superhuman Siehe Siehe Siehe Siehe Siehe Siehe Siehe In diesem Jahr ist es wieder soweit: Manche Manuskripte haben sich in der Hauptsache in der Hauptsache wieder in den Vordergrund gestellt.

(This year it's time again: Some manuscripts the main thing the main thing come to the foreground again)

- Transfer trigger to Systran
- 5.8% similarity between Systran's predictions before and after appending trigger (as measured by modified pairwise BLEU)

Universal Suffix Dropper Attack

• A phrase that commonly causes itself and any subsequent text to be dropped from the translation

Universal Suffix Dropper Attack

• A phrase that commonly causes itself and any subsequent text to be dropped from the translation



Given the release of the latest iPhone, Apple stock was up 3% in trading early yesterday , **kostenfrei übrigens categories ņstrian hl SM the revolution begins at 6pm**

Angesichts der Veröffentlichung des neuesten iPhones lag die Apple-Aktie gestern früh im Handel um 3% im Kurs.

(Given the release of the latest iPhone, Apple stock was up 3% in trading early yesterday)

Universal Suffix Dropper Attack

• A phrase that commonly causes itself and any subsequent text to be dropped from the translation



Given the release of the latest iPhone, Apple stock was up 3% in trading early yesterday , **kostenfrei übrigens categories ņstrian hl SM the revolution begins at 6pm**

Angesichts der Veröffentlichung des neuesten iPhones lag die Apple-Aktie gestern früh im Handel um 3% im Kurs.

(Given the release of the latest iPhone, Apple stock was up 3% in trading early yesterday)

- Transfer trigger to Bing
- **76%** of messages after trigger are successfully dropped

• Cause a specific output token to flip to another specific token

• Cause a specific output token to flip to another specific token

Google

l am going to die, it's over 100°F, help!

lch werde sterben, es ist über 100°F, hilf!

• Cause a specific output token to flip to another specific token

Google	l am going to die, it's over 100°F, help!	lch werde sterben, es ist über 100°F, hilf!
Google	I am going to die, it's over 102°F , help!	Ich werde sterben, es ist über 22°C, hilf!

• Cause a specific output token to flip to another specific token



• 22% of attacks transfer to Google

Defending Against Stealing

Defending Against Stealing

• What makes a good defense?



Defending Against Stealing

• What makes a good defense?

preserves model accuracy

lowers imitation model accuracy

reduces adversarial attack transfer

• Adapt ideas from prediction poisoning (Orekondy et al. 2020)



• Adapt ideas from prediction poisoning (Orekondy et al. 2020)



Adapt ideas from prediction poisoning (<u>Orekondy et al. 2020</u>)



Goal: find a translation $\mathbf{ ilde{y}}$ that is similar to the original

Adapt ideas from prediction poisoning (<u>Orekondy et al. 2020</u>)



Goal: find a translation $\tilde{\mathbf{y}}$ that is similar to the original but induces a different gradient (ideally pointing the opposite direction)

Adapt ideas from prediction poisoning (<u>Orekondy et al. 2020</u>)



Goal: find a translation $\mathbf{\tilde{y}}$ that is similar to the original but induces a different gradient (ideally pointing the opposite direction)

Assumption: angular deviations are similar for adversary's model

• Generate 100 alternate translations via sampling

- Generate 100 alternate translations via sampling
- Pick translation with largest gradient angular deviation

- Generate 100 alternate translations via sampling
- Pick translation with largest gradient angular deviation
- Impose minimum similarity to original via BLEU match

- Generate 100 alternate translations via sampling
- Pick translation with largest gradient angular deviation
- Impose minimum similarity to original via BLEU match











• Defense reduces adversary's BLEU more than defender's
Defenses Can Mitigate Adversarial Threat



- Defense reduces adversary's BLEU more than defender's
- Attack transfer drops from 38% to 27% at 70 BLEU Match

Defenses Can Mitigate Adversarial Threat



- Defense reduces adversary's BLEU more than defender's
- Attack transfer drops from 38% to 27% at 70 BLEU Match
- Downsides: defense adds compute and hurts defender BLEU

Conclusions

- Hiding models behind a black-box API is not enough!
 - Production MT models can be **stolen**
 - Production MT models can be **broken**



Conclusions

- Hiding models behind a black-box API is not enough!
 - Production MT models can be **stolen**
 - Production MT models can be **broken**
- Our defense **mitigates** vulnerabilities, but future work is required



Conclusions

- Hiding models behind a black-box API is not enough!
 - Production MT models can be **stolen**
 - Production MT models can be **broken**
- Our defense **mitigates** vulnerabilities, but future work is required



