

# Automatic Information Extraction from Unstructured Text – *practices and case studies*

---

SHUBHANSHU MISHRA

ISCHOOL AT ILLINOIS

@THESHUBHANSHU

# Structured v/s Unstructured Text

## Structured Text

**Persons:** Rafa → wiki(Rafael\_Nadal), Fed → wiki(Roger Federer)

**Event:** USOpen → wiki(US Open 2017)

**Time:** May 2<sup>nd</sup>, 2017

**Sentiment:** Support towards Rafa

- Rarely available
- Unique representation
- Easy to parse and load in database
- Easy to answer questions on
- Easy to search

## Un-Structured Text

**Rafa and Fed take on each other on May 2<sup>nd</sup>**  
**#USOpen #GoRafa.**

*Timestamp: April 31<sup>st</sup>, 2017*

- Abundantly available
- Multiple representations
- Easy to parse if data is proper English (e.g. newswire corpus)
- Harder to parse when the language used is different

# Information extraction issues

Source: <https://twitter.com/BillSimons1/status/854713304949825540>

— Text to annotate —

The 10s Gods Must B Crazy/Could we have a10s wrld where #Federer, 35, wins 2 huge tourneys & King of Clay, Senor #Nadal, loses in 2nd rd.

— Annotations —

parts-of-speech x named entities x wikipedia entities x sentiment x

— Language —

English

Submit

## Part-of-Speech:

Source: <http://corenlp.run/>

DT CD NNS MD NN NN PRP VBP CD NN WRB NNP CD VBZ CD JJ NNS CC NNP IN NNP NNP NNP VBZ IN JJ NN  
1 The 10s Gods Must B Crazy/Could we have a10s wrld where #Federer , 35 , wins 2 huge tourneys & King of Clay , Senor #Nadal , loses in 2nd rd. .

## Named Entity Recognition:

NUMBER Tennis? NUMBER Federer NUMBER 35 NUMBER 2 Nadal ORDINAL 2.0  
1 The 10s Gods Must B Crazy/Could we have a10s wrld where #Federer , 35 , wins 2 huge tourneys & King of Clay , Senor #Nadal , loses in 2nd rd. .

## Wikidict Entities:

10s 35 2  
1 The 10s Gods Must B Crazy/Could we have a10s wrld where #Federer , 35 , wins 2 huge tourneys & King of Clay , Senor #Nadal , loses in 2nd rd. .

## Sentiment:

NEGATIVE  
1 The 10s Gods Must B Crazy/Could we have a10s wrld where #Federer , 35 , wins 2 huge tourneys & King of Clay , Senor #Nadal , loses in 2nd rd. .

- Is it even right to assign a polarity to the overall text?
- What does positive or negative sentiment mean? (taken from ML lingo)
- Shouldn't sentiment convey action or internal state of—either the reader or the author?

# Traditional NLP tools break on Social Media Data

— Text to annotate —  
Feds is mah man :) Rafa is a loser #burn #USOpen

— Annotations —  
parts-of-speech x named entities x sentiment x wikipedia entities x

Part-of-Speech: **Emoticon lost**

1	Feds	is	mah	man	:	-RRB-	Rafa	is	a	loser	#burn	#USOpen
	NNS	VBZ	JJ	NN	NN	NN	VBZ	DT	JJR	NN	NN	

Named Entity Recognition:

1	Feds is mah man :-RRB- Rafa is a loser #burn #USOpen
---	------------------------------------------------------

Wikidict Entities:

1	Feds is mah man :-RRB- Rafa is a loser #burn #USOpen
---	------------------------------------------------------

Sentiment:

1	Feds is mah man :-RRB- Rafa is a loser #burn #USOpen
	NEGATIVE

Source: <http://corenlp.run/>

— Text to annotate —  
Feds is my man. Rafa is a loser.

Easier to parse and tag

— Annotations —  
parts-of-speech x named entities x sentiment x wikipedia entities x

Part-of-Speech:

1	Feds	is	my	man	.
	NNS	VBZ	PRPS	NN	.
2	Rafa	is	a	loser	.
	NNP	VBZ	DT	JJR	.

Named Entity Recognition:

1	Feds is my man .
2	Rafa is a loser .

Wikidict Entities:

1	Feds is my man .
2	Rafa is a loser .

Sentiment:

1	Feds is my man .
	NEUTRAL
2	Rafa is a loser .
	NEUTRAL

**Sentiment lost**

*“Our field is the **domain science of language technology**; it’s **not** about the best method of machine learning—the central issue remains the **domain problems**. The domain problems **will not go away**.”*

**--- Christopher Manning** (ACL President 2015)

Source: <http://mitp.nautil.us/article/170/last-words-computational-linguistics-and-deep-learning>

# Solutions for NLP in Social Media Text

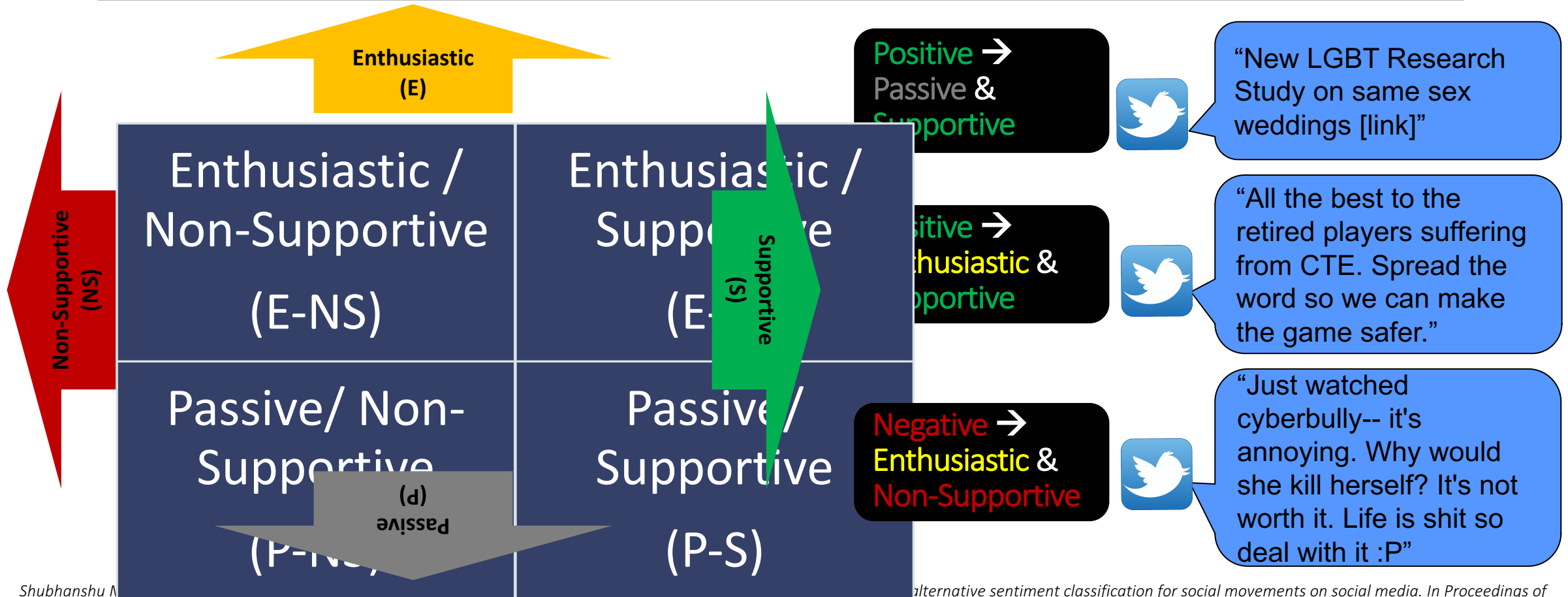
## *Sentiment Analysis:*

- Define **actionable labels** for sentiment to capture if the author is:
  - **Enthusiastic or passive** in their expression
  - **Supportive or non-supportive** of the topic
- Make task specific labels – don't just stick to positive or negative
- Utilize the **context of the text** for making prediction:
  - Who is the **user**? – number of followers, friends, statuses, etc.
  - What **metadata** is associated with the tweet? – retweets, favorites, etc.

## *Named Entity Recognition and Part of Speech Tagging:*

- Allow models to:
  - Change with **time**
  - Adapt to your testing data
- This can be achieved via:
  - Updating the **vocabulary** of the model
  - Using **online or semi-supervised learning** models to learn from unlabeled and newer data

# Actionable labels: *Enthusiasm and Support*



# Actionable labels:

## *Tuition-Free Higher Education Policy Debate on Facebook*

Opinion on subject

- **For:** Tax money used for a good long-term investment, finally.
- **Against:** Free? It's not free-someone's gotta pay.

Tone of the conversation

- **Civil:** Nothing is free Mr. President.
- **Non-Civil:** NOTHING IS FREE!!!!!!! WAKE UP!!

Relatedness to the topic

- **On-Topic:** I hope they pay for it with cuts in wasteful spending, not more taxes.
- **Off-Topic:** There's ALWAYS a catch, just like with "You can keep your insurance"



# Use context

*User attributes: Gender, Ethnicity, etc.*

---

## Gender

- Extracted via US SSN data
- Can also use existing image based classifiers

## Ethnicity

- Predicted based on full name via Textmap tool. <http://www.textmap.com/ethnicity/>

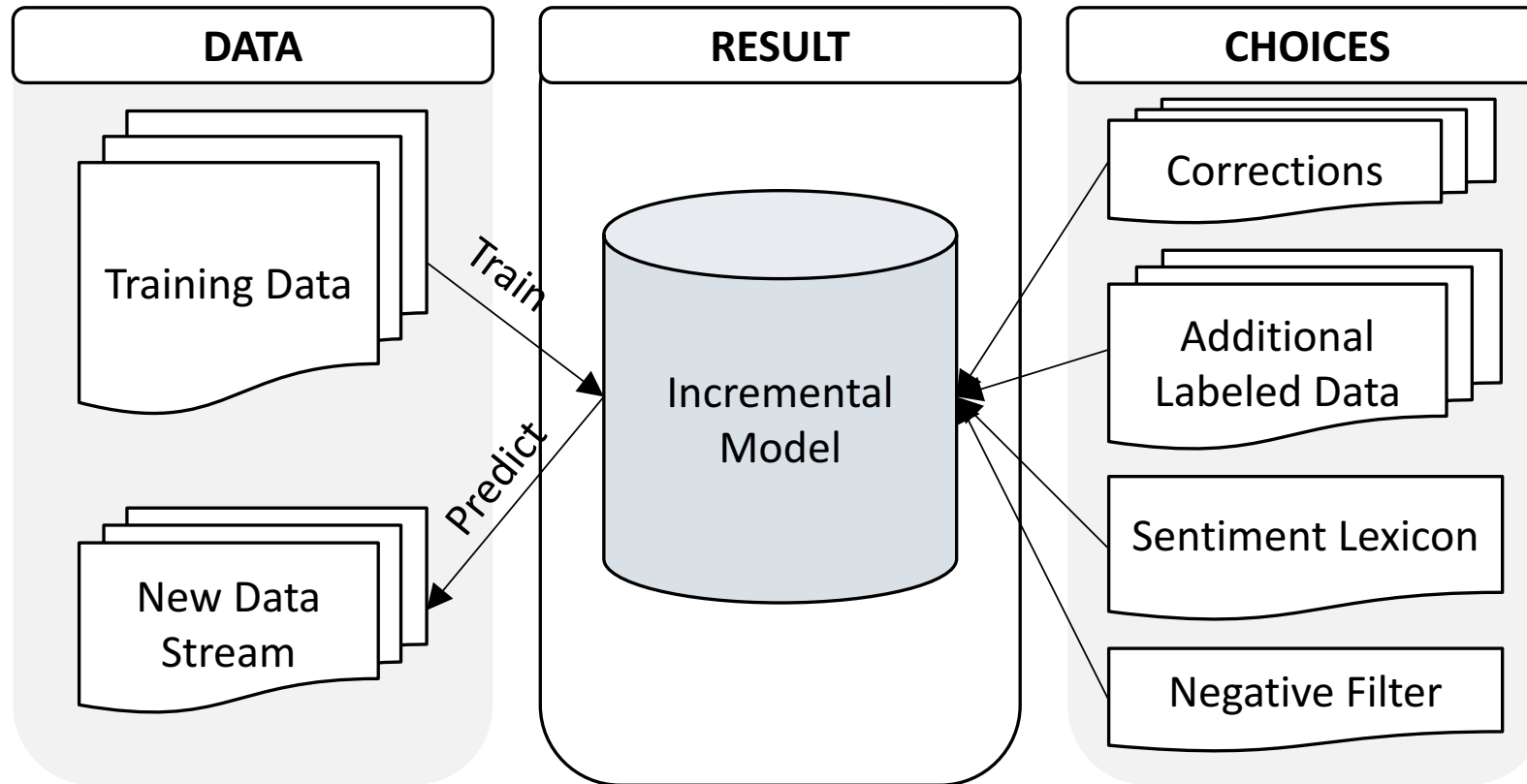
## Race

- Manually annotated using profile inspection

## Political Leaning

- Manually annotated using profile inspection (if available)

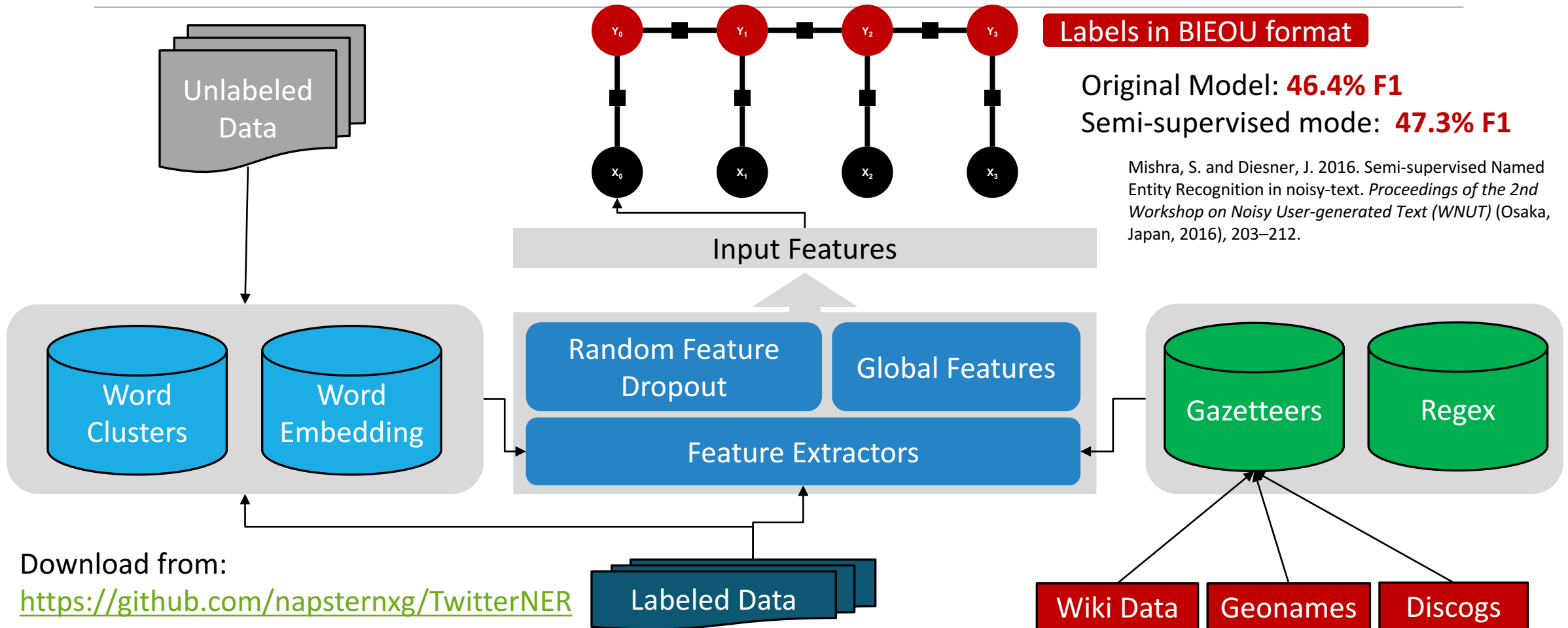
# Adaptive models (Sentiment): *Vocabulary Update & Online Learning*



Download tool from:

<https://github.com/uiuc-ischool-scanr/SAIL>

# Adaptive models (NER): *Vocabulary Update & Semi-supervised Learning*



Download from:

<https://github.com/napsternxg/TwitterNER>

# Open challenges

---

Not enough labeled data for social media IE tasks

Not enough pre-trained models

Rapidly evolving contents

Privacy and TOS issues – data available today might be missing tomorrow

# Key Takeaways

---

Don't apply existing tools blindly

Understand your labels and their relation with the data

Use the context

Allow models to evolve (be aware if they don't)

# Questions ?

---

**Contact:**

Shubhanshu Mishra

<http://shubhanshu.com>

[@TheShubhanshu](#)