# DSC291: Machine Learning with Few Labels

## Weak/distant supervision

**Zhiting Hu**

Lecture 12, February 6, 2023

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

# Recap: Data Augmantation

- Image:
  - *Flip, crop, scale, translation, rotation, mixup, …*

- Text:

| Methods | Level | Diversity | Tasks | Related Work |
|---|---|---|---|---|
| Synonym replacement | Token | Low | Text classification Sequence labeling | Kolomiyets et al. (2011), Zhang et al. (2015a), Yang (2015), Miao et al. (2020), Wei and Zou (2019) |
| Word replacement via LM | Token | Medium | Text classification Sequence labeling Machine translation | Kolomiyets et al. (2011), Gao et al. (2019) Kobayashi (2018), Wu et al. (2019a) Fadaee et al. (2017) |
| Random insertion, deletion, swapping | Token | Low | Text classification Sequence labeling Machine translation Dialogue generation | Iyyer et al. (2015), Xie et al. (2017) Artetxe et al. (2018), Lample et al. (2018) Xie et al. (2020), Wei and Zou (2019) |
| Compositional Augmentation | Token | High | Semantic Parsing Sequence labeling Language modeling Text generation | Jia and Liang (2016) , Andreas (2020) Nye et al. (2020), Feng et al. (2020) Furrer et al. (2020) , Guo et al. (2020) |
| Paraphrasing | Sentence | High | Text classification Machine translation Question answering Dialogue generation Text summarization | Yu et al. (2018), Xie et al. (2020) Chen et al. (2019), He et al. (2020) Chen et al. (2020c), Cai et al. (2020) |
| Conditional generation | Sentence | High | Text classification Question answering | Anaby-Tavor et al. (2020), Kumar et al. (2020) Zhang and Bansal (2019), Yang et al. (2020) |

# Recap: Data Augmantation

- Image:

  - *Flip, crop, scale, translation, rotation, mixup, …*

- Text:

| | | | | |
|---|---|---|---|---|
| White-box attack | Token or Sentence | Medium | Text classification<br>Sequence labeling<br>Machine translation | Miyato et al. (2017), Ebrahimi et al. (2018b)<br>Ebrahimi et al. (2018a), Cheng et al. (2019),<br>Chen et al. (2020d) |
| Black-box attack | Token or Sentence | Medium | Text classification<br>Sequence labeling<br>Machine translation<br>Textual entailment<br>Dialogue generation<br>Text Summarization | Jia and Liang (2017)<br>Belinkov and Bisk (2017), Zhao et al. (2017)<br>Ribeiro et al. (2018), McCoy et al. (2019)<br>Min et al. (2020), Tan et al. (2020) |
| Hidden-space perturbation | Token or Sentence | High | Text classification<br>Sequence labeling<br>Speech recognition | Hsu et al. (2017), Hsu et al. (2018)<br>Wu et al. (2019b), Chen et al. (2021)<br>Malandrakis et al. (2019), Shen et al. (2020) |
| Interpolation | Token | High | Text classification<br>Sequence labeling<br>Machine translation | Miao et al. (2020), Chen et al. (2020c)<br>Cheng et al. (2020b), Chen et al. (2020a)<br>Guo et al. (2020) |

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution
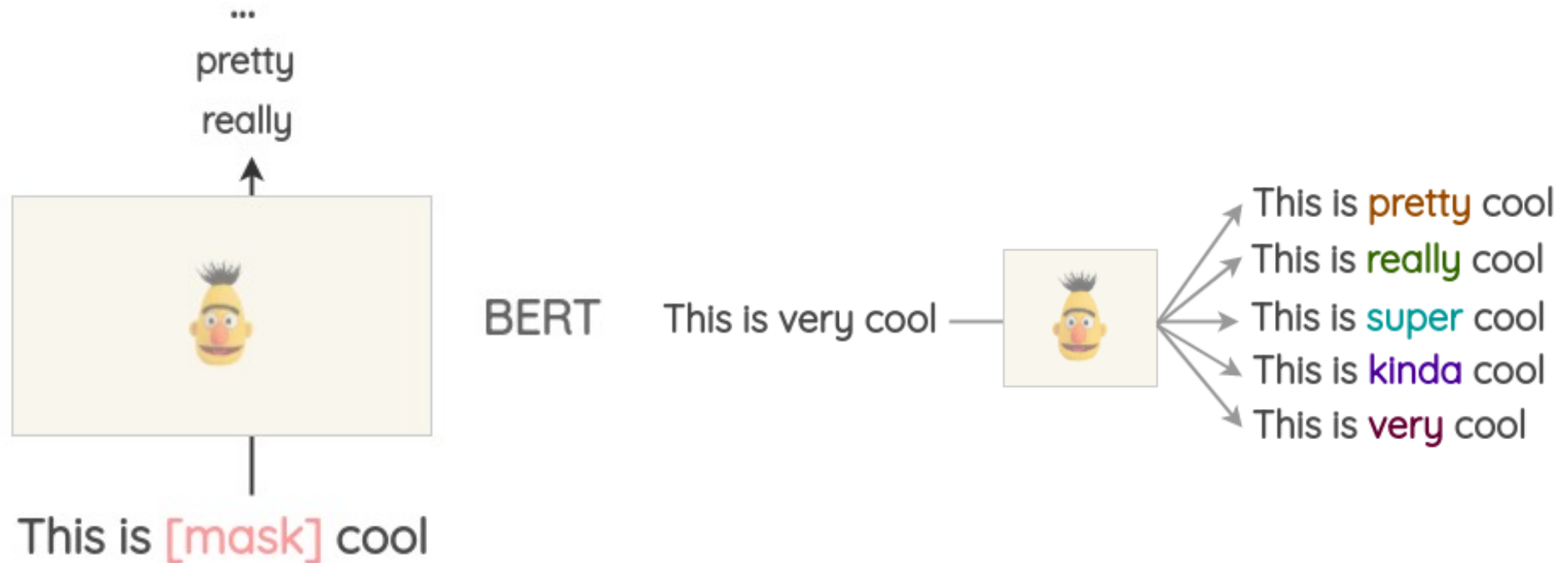  - Word-embedding substitution

Nearest neighbors in word2vec



It is <u>awesome</u>
- It is amazing
- It is perfect
- It is fantastic

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution
  - Word-embedding substitution
  - Masked LM

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution
  - Word-embedding substitution
  - Masked LM
  - TF-IDF based word replacement
    - words that have low TF-IDF scores are uninformative and thus can be replaced without affecting the ground-truth labels of the sentence.
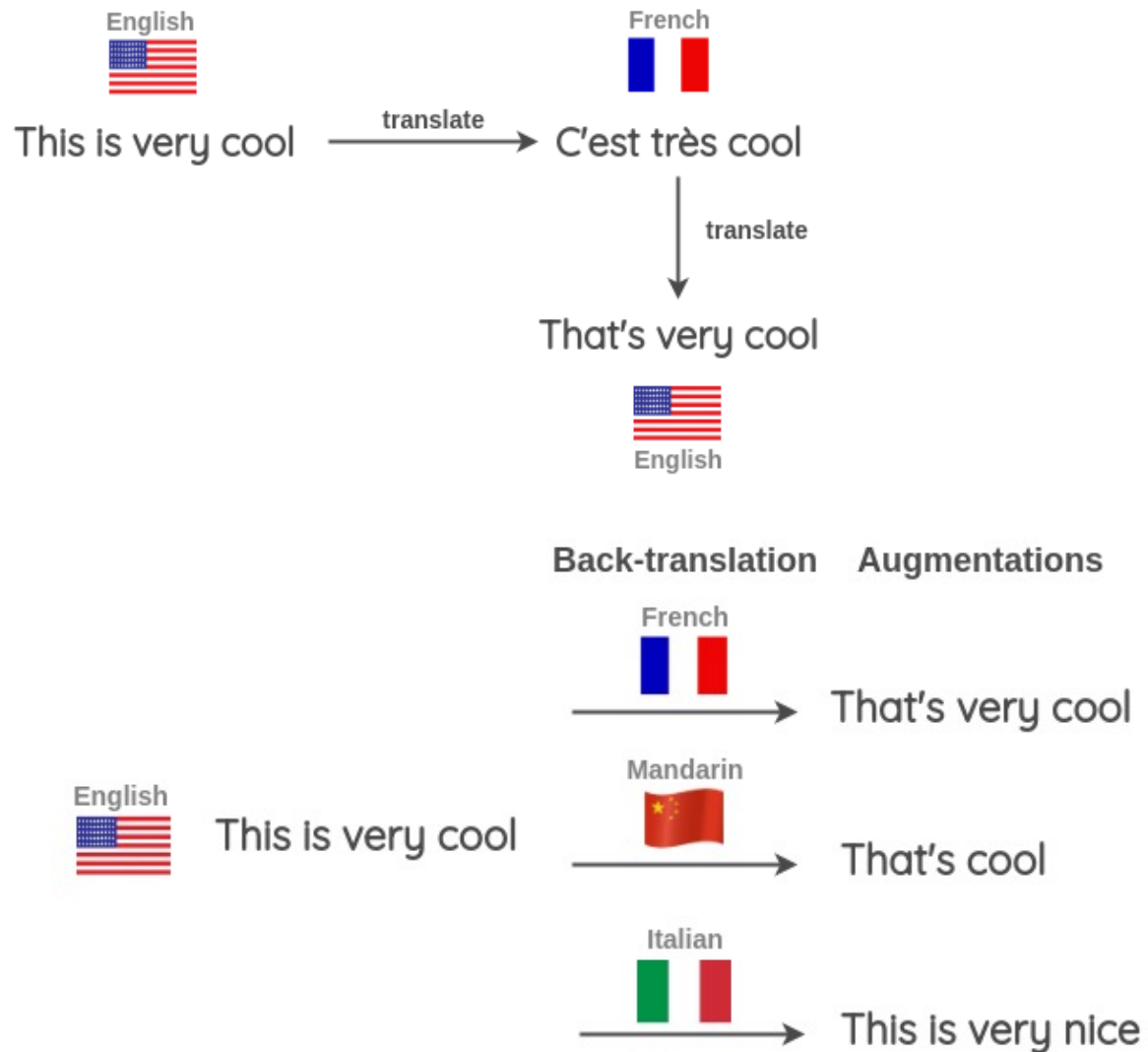
**This** virus has spread worldwide

↓

A virus has spread worldwide

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution
  - Word-embedding substitution
  - Masked LM
  - TF-IDF based word replacement
- Paraphrasing
  - Back Translation

English

This is very cool →translate→ French C'est très cool

translate ↓

That's very cool

English

Back-translation  Augmentations

English This is very cool

French →  That's very cool

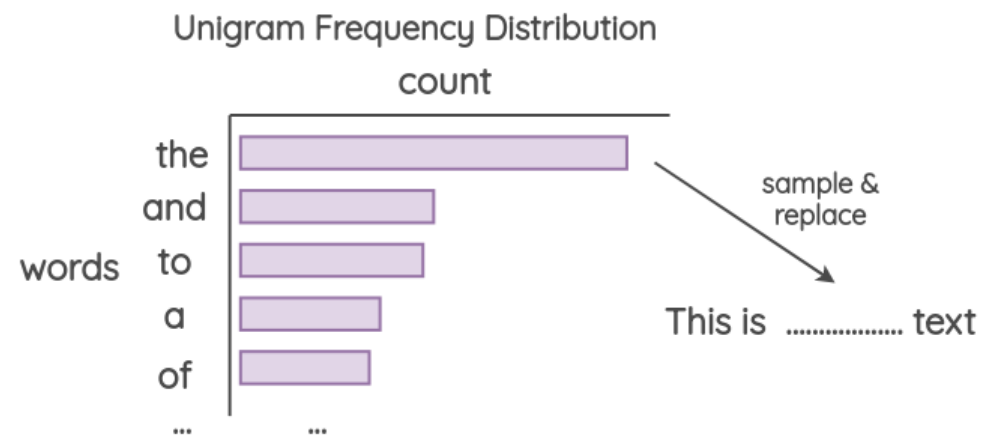Mandarin →  That's cool

Italian →  This is very nice

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution
  - Word-embedding substitution
  - Masked LM
  - TF-IDF based word replacement

- Paraphrasing
  - Back Translation

- Random Noise Injection

Spelling error:

This is very cool → Thes is very cool
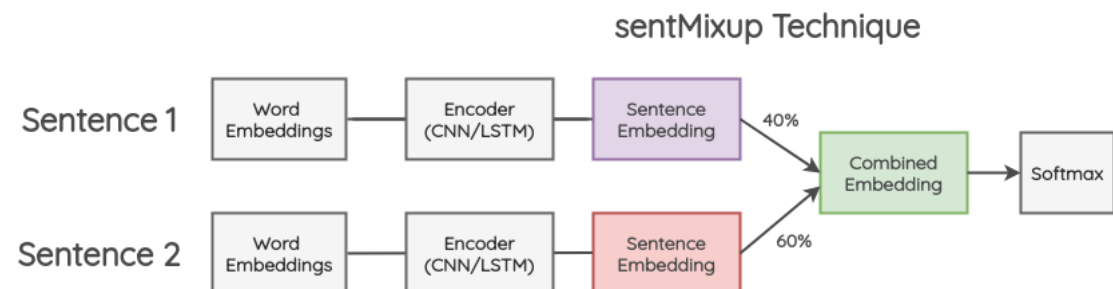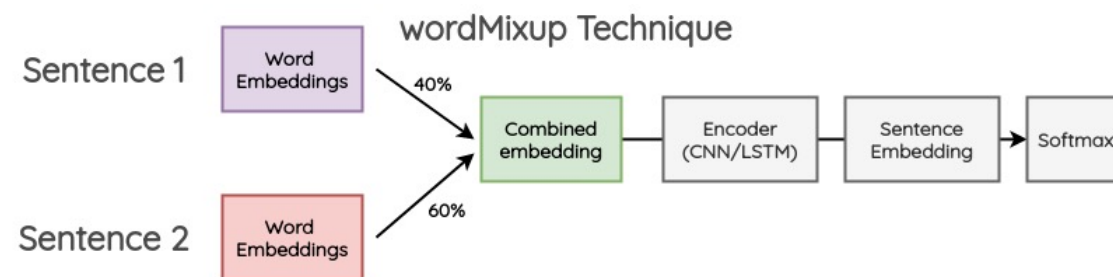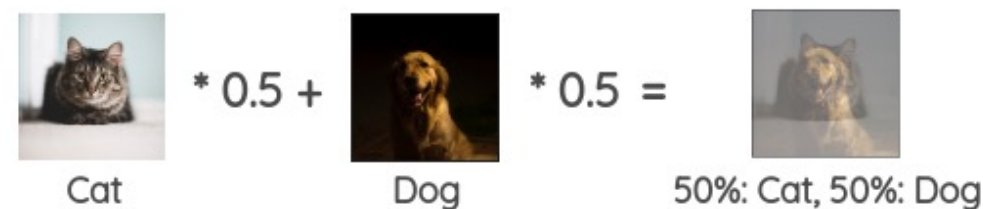
This is very cool → This id very cool

Unigram noising:

Unigram Frequency Distribution

count

words: the, and, to, a, of, ...

sample & replace

This is ............... text

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution
  - Word-embedding substitution
  - Masked LM
  - TF-IDF based word replacement
- Paraphrasing
  - Back Translation
- Random Noise Injection
- MixUp



Original Mixup algorithm

Cat * 0.5 + Dog * 0.5 = 50%: Cat, 50%: Dog

wordMixup Technique

Sentence 1 → Word Embeddings → 40% → Combined embedding → Encoder (CNN/LSTM) → Sentence Embedding → Softmax

Sentence 2 → Word Embeddings → 60%

sentMixup Technique

Sentence 1 → Word Embeddings → Encoder (CNN/LSTM) → Sentence Embedding → 40% → Combined Embedding → Softmax

Sentence 2 → Word Embeddings → Encoder (CNN/LSTM) → Sentence Embedding → 60%

# Text data augmentation: Examples

- Lexical Substitution
  - Thesaurus-based substitution
  - Word-embedding substitution
  - Masked LM
  - TF-IDF based word replacement
- Paraphrasing
  - Back Translation
- Random Noise Injection
- MixUp
- Generative Models
  - Finetune a large pre-trained LM (BERT, GPT2, etc)
  - Use the fine-tuned LM to generate new data

Finetune on training data

GPT2

Task: Learn to generate training data
Output: POSITIVE<SEP>It is very useful app<EOS>

Generate new samples
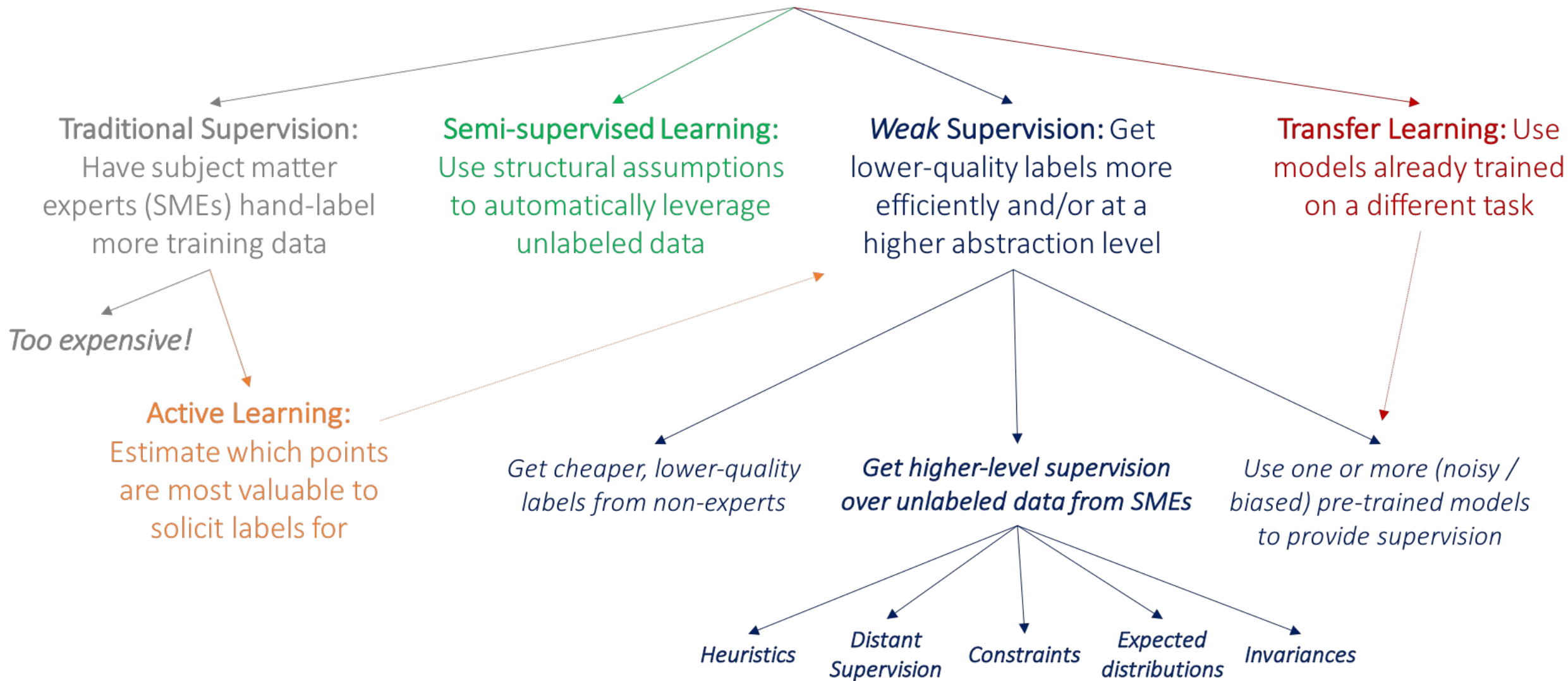
GPT2

Prompt: POSITIVE <SEP>It is very
Generate: POSITIVE <SEP> It is very helpful tool<EOS>

# Weakly Supervised Learning

# The difficulty with supervised learning

- Annotated data is expensive and costs increase when…

  - *A task requires specialized expertise*

    *E.g. "Only a trained linguist or a board certified radiologist can label my data"*

  - *Labeling examples involves making multiple decisions*

    *E.g. "Annotate this sentence with a parse tree"*

    *(instead of a single binary decision)*

# How to get more labeled training data?



**Traditional Supervision:** Have subject matter experts (SMEs) hand-label more training data

*Too expensive!*

**Semi-supervised Learning:** Use structural assumptions to automatically leverage unlabeled data

***Weak* Supervision:** Get lower-quality labels more efficiently and/or at a higher abstraction level

**Transfer Learning:** Use models already trained on a different task

**Active Learning:** Estimate which points are most valuable to solicit labels for

*Get cheaper, lower-quality labels from non-experts*

***Get higher-level supervision over unlabeled data from SMEs***

*Use one or more (noisy / biased) pre-trained models to provide supervision*

*Heuristics*   *Distant Supervision*   *Constraints*   *Expected distributions*   *Invariances*

# Example (I): labeling with heuristics

Task: Build a chest x-ray classifier
(normal/abnormal)



Indication: Chest pain. Findings:
Mediastinal contours are within
`normal` limits. Heart size is
within `normal` limits. `No` focal
consolidation, `pneumothorax` or
`pleural effusion`. Impression: `No`
acute cardiopulmonary
abnormality.

Can you use the accompanying medical report (text modality) to label the x-ray (image modality)?

# Example (I): labeling with heuristics

Indication: Chest pain. Findings: Mediastinal contours are within `normal` limits. Heart size is within `normal` limits. `No` focal consolidation, `pneumothorax` or `pleural effusion`. Impression: `No` acute cardiopulmonary abnormality.

How do we obtain Y?

Y

CNN

# Example (I): labeling with heuristics

Indication: Chest pain. Findings:
Mediastinal contours are within
`normal` limits. Heart size is
within `normal` limits. `No` focal
consolidation, `pneumothorax` or
`pleural effusion`. Impression: `NO`
acute cardiopulmonary
abnormality.

**Normal Report**

```python
def LF_pneumothorax(c):
    if re.search(r'pneumo.*', c.report.text):
        return "ABNORMAL"

def LF_pleural_effusion(c):
    if "pleural effusion" in c.report.text:
        return "ABNORMAL"

def LF_normal_report(c, thresh=2):
    if len(NORMAL_TERMS.intersection(c.
    report.words)) > thresh:
        return "NORMAL"
```

**LFs**

(labeling functions)

Source: Khandwala et. al 2017, Cross Modal Data Programming for Medical Images

# Example (II): Labeling with knowledge bases

Task: relation extraction from text

- Hypothesis: If two entities belong to a certain relation, any sentence containing those two entities is likely to express that relation
- Key idea: use a *knowledge base* of relations to get lots of *noisy* training examples

# Example (II): Labeling with knowledge bases
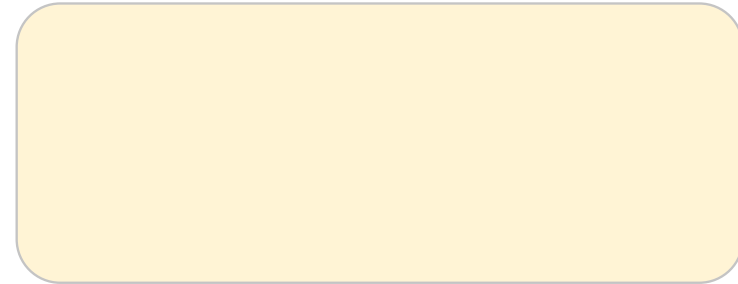
## Frequent Freebase relations

| Relation name | Size | Example |
|---|---|---|
| /people/person/nationality | 281,107 | John Dugard, South Africa |
| /location/location/contains | 253,223 | Belgium, Nijlen |
| /people/person/profession | 208,888 | Dusa McDuff, Mathematician |
| /people/person/place_of_birth | 105,799 | Edwin Hubble, Marshfield |
| /dining/restaurant/cuisine | 86,213 | MacAyo's Mexican Kitchen, Mexican |
| /business/business_chain/location | 66,529 | Apple Inc., Apple Inc., South Park, NC |
| /biology/organism_classification_rank | 42,806 | Scorpaeniformes, Order |
| /film/film/genre | 40,658 | Where the Sidewalk Ends, Film noir |
| /film/film/language | 31,103 | Enter the Phoenix, Cantonese |
| /biology/organism_higher_classification | 30,052 | Calopteryx, Calopterygidae |
| /film/film/country | 27,217 | Turtle Diary, United States |
| /film/writer/film | 23,856 | Irving Shulman, Rebel Without a Cause |
| /film/director/film | 23,539 | Michael Mann, Collateral |
| /film/producer/film | 22,079 | Diane Eskenazi, Aladdin |
| /people/deceased_person/place_of_death | 18,814 | John W. Kern, Asheville |
| /music/artist/origin | 18,619 | The Octopus Project, Austin |
| /people/person/religion | 17,582 | Joseph Chartrand, Catholicism |
| /book/author/works_written | 17,278 | Paul Auster, Travels in the Scriptorium |
| /soccer/football_position/players | 17,244 | Midfielder, Chen Tao |
| /people/deceased_person/cause_of_death | 16,709 | Richard Daintree, Tuberculosis |
| /book/book/genre | 16,431 | Pony Soldiers, Science fiction |
| /film/film/music | 14,070 | Stavisky, Stephen Sondheim |
| /business/company/industry | 13,805 | ATS Medical, Health care |

# Example (II): Labeling with knowledge bases

## Corpus text

Bill Gates founded Microsoft in 1975.

Bill Gates, founder of Microsoft, …

Bill Gates attended Harvard from…
Google was founded by Larry Page …

## Training data

## Freebase

Founder: (Bill Gates, Microsoft)

Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

# Example (II): Labeling with knowledge bases

## Corpus text

Bill Gates founded Microsoft in 1975.

Bill Gates, founder of Microsoft, …

Bill Gates attended Harvard from…

Google was founded by Larry Page …

## Training data

(Bill Gates, Microsoft)
Label:         Founder
Feature:     X founded Y

## Freebase

Founder: (Bill Gates, Microsoft)

Founder: (Larry Page, Google)

CollegeAttended: (Bill Gates, Harvard)

# Example (II): Labeling with knowledge bases

## Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, …
Bill Gates attended Harvard from…
Google was founded by Larry Page …

## Training data

(Bill Gates, Microsoft)
Label:        Founder
Feature:    X founded Y
Feature:    X, founder of Y

## Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

# Example (II): Labeling with knowledge bases

## Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, …
Bill Gates attended Harvard from…
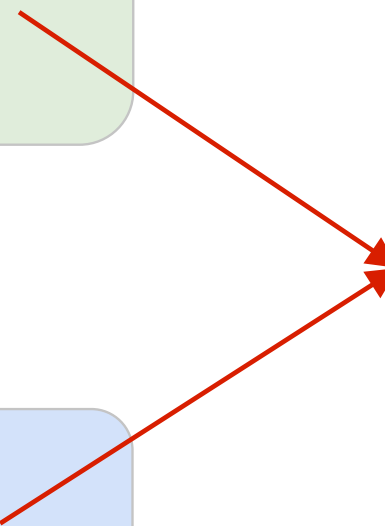Google was founded by Larry Page …

## Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

## Training data

(Bill Gates, Microsoft)
Label:       Founder
Feature:     X founded Y
Feature:     X, founder of Y

(Bill Gates, Harvard)
Label:       CollegeAttended
Feature:     X attended Y

# Example (II): Labeling with knowledge bases

## Corpus text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, …
Bill Gates attended Harvard from…
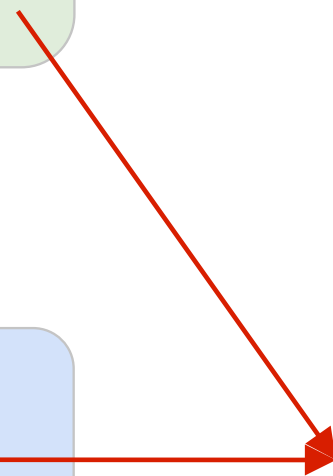Google was founded by Larry Page …

## Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

## Training data

(Bill Gates, Microsoft)
Label:       Founder
Feature:     X founded Y
Feature:     X, founder of Y

(Bill Gates, Harvard)
Label:       CollegeAttended
Feature:     X attended Y

(Larry Page, Google)
Label:       Founder
Feature:     Y was founded by X

# Example (II): Labeling with knowledge bases

## Negative training data

Can't train a classifier with only positive data!
Need negative training data too!

Solution?
Sample 1% of unrelated pairs of entities.
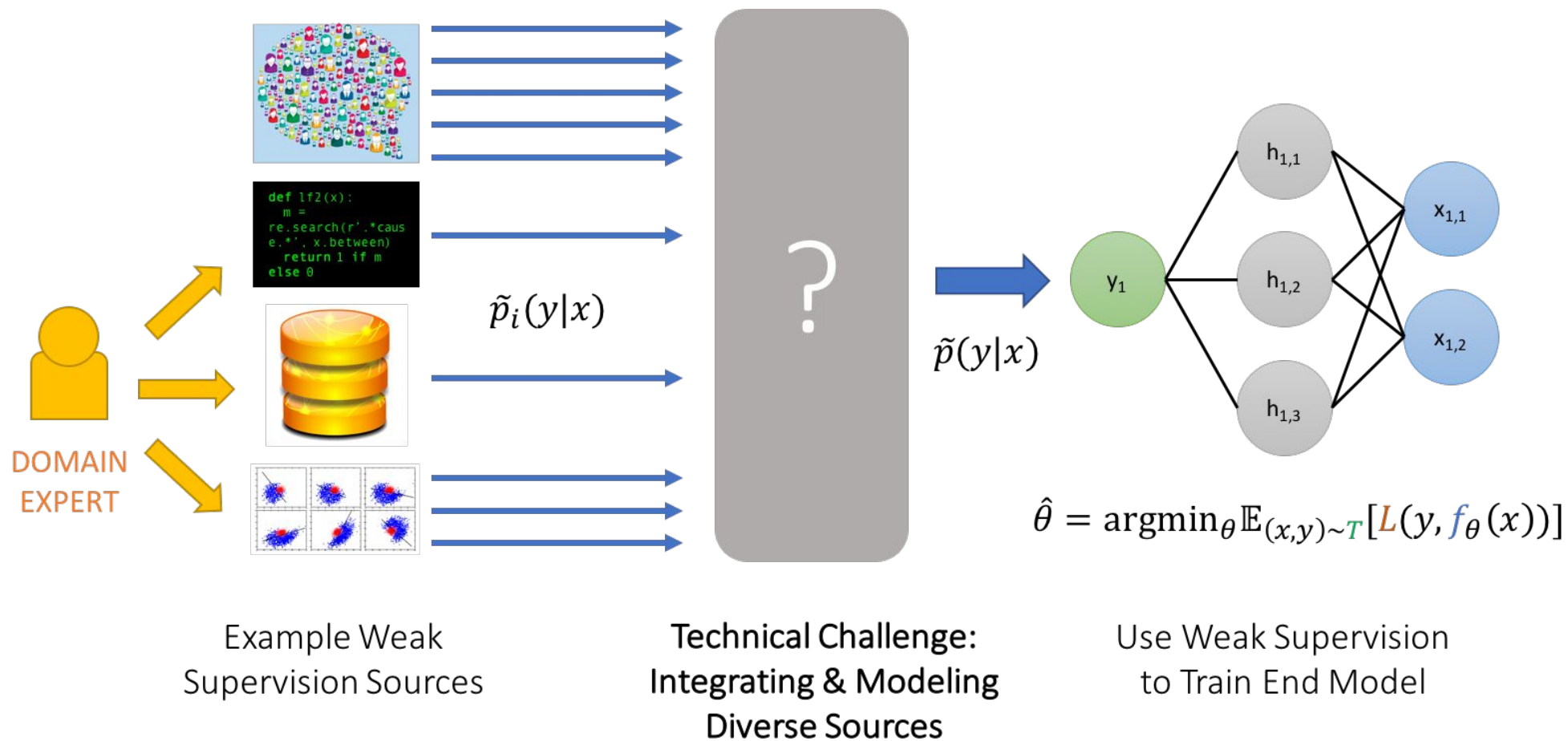
### Training data

(Larry Page, Microsoft)
Label:       NO_RELATION
Feature:     X took a swipe at Y

(Larry Page, Harvard)
Label:       NO_RELATION
Feature:     Y invited X

(Bill Gates, Google)
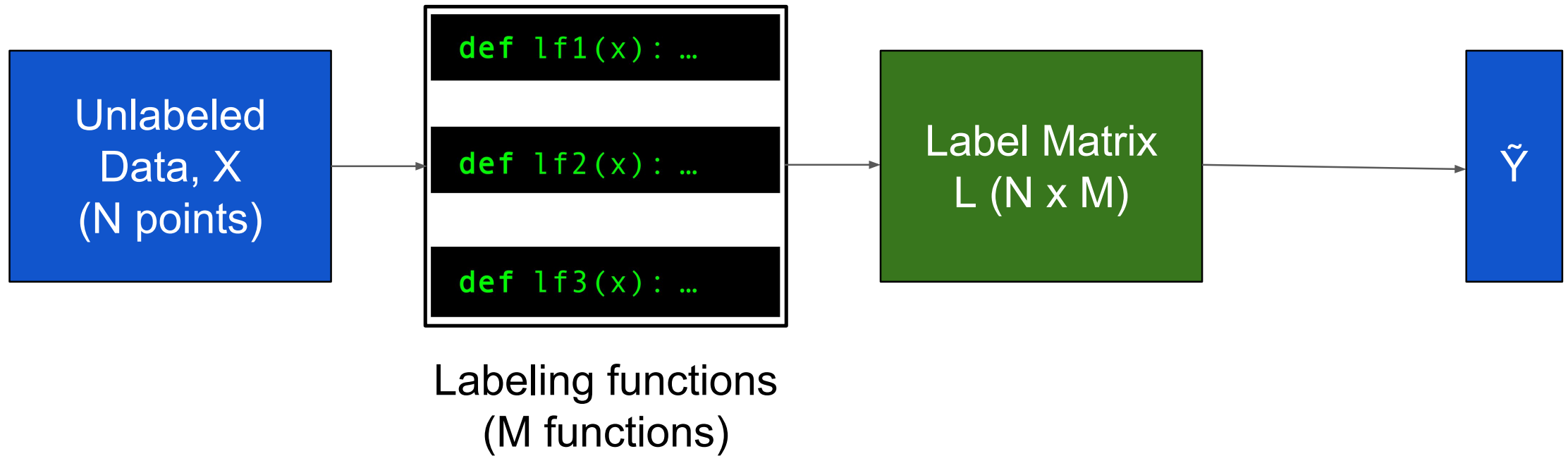Label:       NO_RELATION
Feature:     Y is X's worst fear

### Corpus text

Larry Page took a swipe at Microsoft...
...after Harvard invited Larry Page to...
Google is Bill Gates' worst fear ...

# Integrating multiple noisy labels



$\tilde{p}_i(y|x)$

$\tilde{p}(y|x)$

$$\hat{\theta} = \text{argmin}_\theta \mathbb{E}_{(x,y)\sim T}[L(y, f_\theta(x))]$$

Example Weak
Supervision Sources

Technical Challenge:
Integrating & Modeling
Diverse Sources

Use Weak Supervision
to Train End Model

# Integrating multiple noisy labels

# Integrating multiple noisy labels



Labeling functions
(M functions)

# Integrating multiple noisy labels

How do we obtain probabilistic labels, **Ỹ,** from the label matrix, **L**?

Approach 1 - Majority Vote

Take the majority vote of the labelling functions (LFs).

Let's say **L** = [[0, 1, 0, 1, 0]; [1, 1, 1, 1, 0]].

$$\tilde{\textbf{Y}} = [0, 1]$$

# Integrating multiple noisy labels

How do we obtain probabilistic labels, $\tilde{Y}$, from the label matrix, $L$?

## Approach 1 - Majority Vote

*Majority vote fails:*

Indication: Chest pain. Findings:
Mediastinal contours are within
`normal` limits. Heart size is
within `normal` limits. `No` focal
consolidation, `pneumothorax` or
`pleural effusion`. Impression: `NO`
acute cardiopulmonary
abnormality.

**Normal Report**

```python
def LF_pneumothorax(c):
    if re.search(r'pneumo.*', c.report.text):
        return "ABNORMAL"

def LF_pleural_effusion(c):
    if "pleural effusion" in c.report.text:
        return "ABNORMAL"

def LF_normal_report(c, thresh=2):
    if len(NORMAL_TERMS.intersection(c.
    report.words)) > thresh:
        return "NORMAL"
```
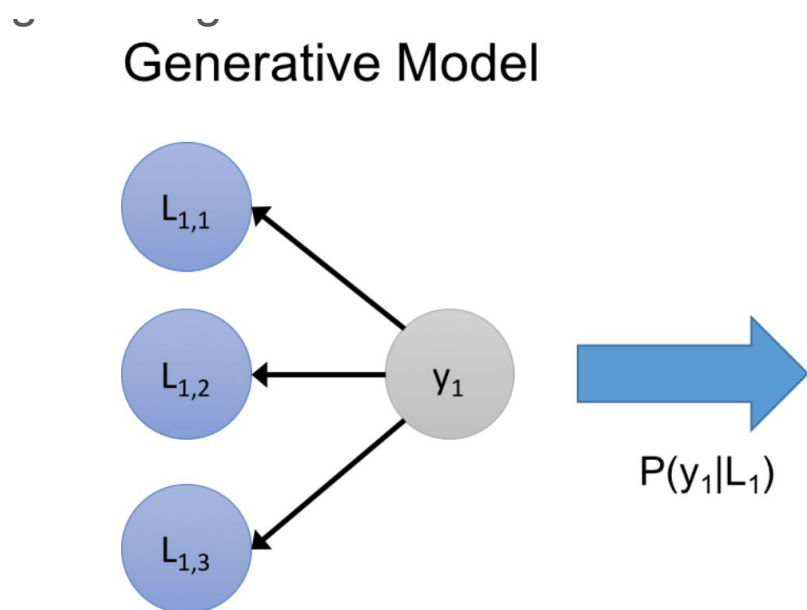
**LFs**

# Integrating multiple noisy labels

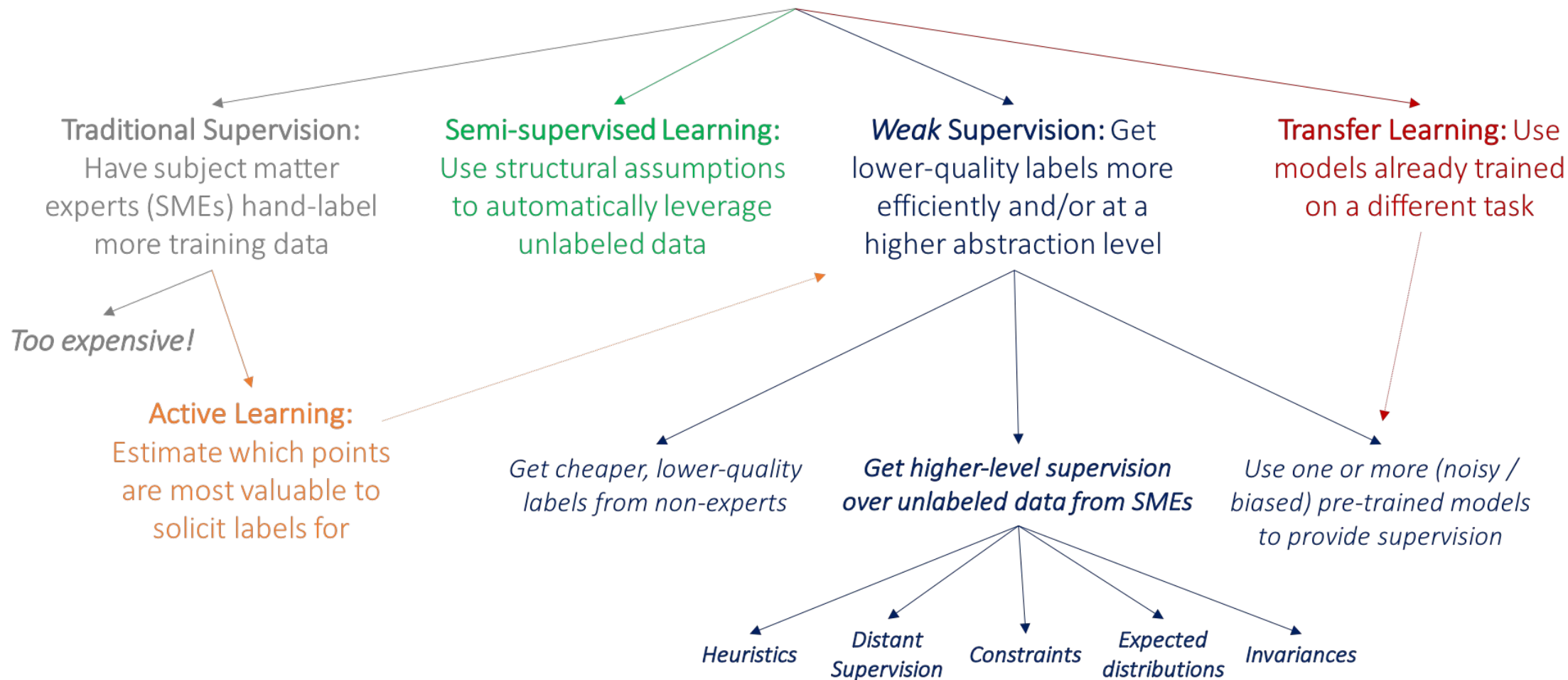How do we obtain probabilistic labels, $\tilde{Y}$, from the label matrix, $L$?

## Approach 2

Train a generative model over $P(L, Y)$ where $Y$ are the **(unknown)** true labels



Generative Model

$P(y_1|L_1)$

# Summary: Weak/distant supervision

How to get more labeled training data?

**Traditional Supervision:** Have subject matter experts (SMEs) hand-label more training data

*Too expensive!*

**Active Learning:** Estimate which points are most valuable to solicit labels for

**Semi-supervised Learning:** Use structural assumptions to automatically leverage unlabeled data

***Weak* Supervision:** Get lower-quality labels more efficiently and/or at a higher abstraction level

**Transfer Learning:** Use models already trained on a different task

*Get cheaper, lower-quality labels from non-experts*

*Get higher-level supervision over unlabeled data from SMEs*

*Use one or more (noisy / biased) pre-trained models to provide supervision*

*Heuristics*     *Distant Supervision*     *Constraints*     *Expected distributions*     *Invariances*

# Summary: Weak/distant supervision

- Noisy labels from heuristics, knowledge bases, constraints, …
- Integrating multiple noisy labels
  - Majority vote
  - Generative modeling
  - …


- Not all information/experiences can easily be converted into labels
  - "Every part of speech sequence should have a verb"
  - "In a sentence with word 'but', the sentiment of text after 'but' dominates"
  - "Every image patch that is recognized as a bicycle should have at least one patch that is recognized as a wheel"
  - I have a "discriminator" model that can tell me whether a model-generated image is good or not
- Need a more flexible framework to incorporate all forms of experiences

# Questions?