# **Fonduer**: Knowledge Base Construction from Richly Formatted Data

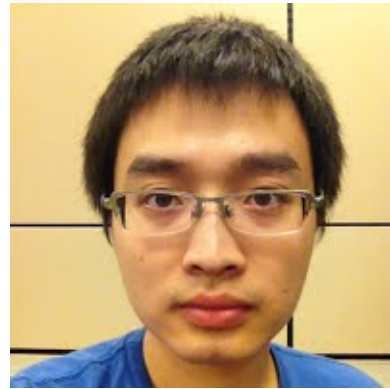## **Sen Wu**

Stanford University

# Thank you!

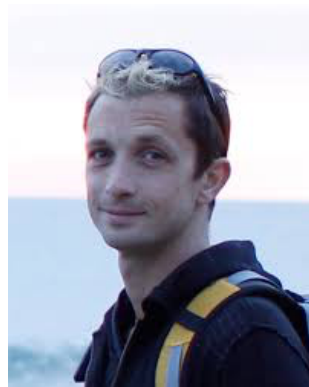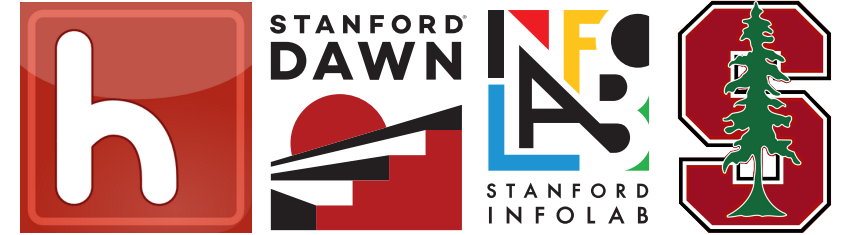FONDUER

## My Amazing Team:

Luke Hsiao

Xiao Cheng

Braden Hancock

Theodoros Rekatsinas

Philip Levis

Christopher Ré

## My Awesome Collaborators:

# Knowledge bases are everywhere...

**Unstructured Information**

Knowledge Base Construction →

**Structured Knowledge Base**

KNOWLEDGE VAULT

yago select knowledge

ProbKB

IBM Watson

Product Graph

TextRunner/ ReVerb

NELL

DeepDive

And many more...

**But, troves of "richly formatted" information remains untapped**

# Richly formatted data

**Richly formatted data**: information is expressed via textual, structural, tabular, and visual cues.

HTML   XML   PDF   DOC

**Transistor Datasheet (PDF)**

**SMBT3904...MMBT3904**

**NPN Silicon Switching Transistors**
- High DC current gain: 0.1 mA to 100 mA
- Low collector-emitter saturation voltage

**Maximum Ratings**

| Parameter | Symbol | Value | Unit |
|---|---|---|---|
| Collector-emitter voltage | $V_{CEO}$ | 40 | V |
| Collector-base voltage | $V_{CBO}$ | 60 | |
| Emitter-base voltage | $V_{EBO}$ | 6 | |
| Collector current | $I_C$ | 200 | mA |
| Total power dissipation | $P_{tot}$ | | mV |
| $T_S \le 71°C$ | | 330 | |
| $T_S \le 115°C$ | | 250 | |
| Junction temperature | $T_i$ | 150 | °C |
| Storage temperature | $T_{stg}$ | -65 ... 150 | |

## Goal: extract maximum collector current from transistor datasheets
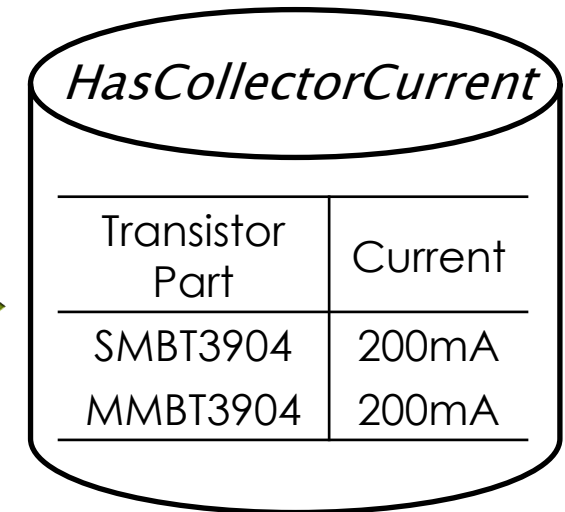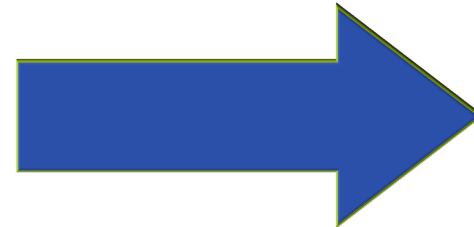
**Transistor Datasheet**

**SMBT3904** ... **MMBT3904**

**NPN Silicon Switching Transistors**
- High DC current gain: 0.1 mA to 100 mA
- Low collector-emitter saturation voltage

**Maximum Ratings**

| Parameter | Symbol | Value | Unit |
|---|---|---|---|
| Collector-emitter voltage | $V_{CEO}$ | 40 | V |
| Collector-base voltage | $V_{CBO}$ | 60 | |
| Emitter-base voltage | $V_{EBO}$ | 6 | |
| Collector current | $I_C$ | 200 | mA |
| Total power dissipation | $P_{tot}$ | | mV |
| $T_S \leq 71°C$ | | 330 | |
| $T_S \leq 115°C$ | | 250 | |
| Junction temperature | $T_i$ | 150 | °C |
| Storage temperature | $T_{stg}$ | -65 ... 150 | |

**HasCollectorCurrent**

| Transistor Part | Current |
|---|---|
| SMBT3904 | 200mA |
| MMBT3904 | 200mA |

Knowledge Base

# Knowledge base construction from richly formatted data



**Transistor Datasheet**

In richly formatted data, **semantics are expressed in** textual, structural, tabular, **and** visual **modalities throughout a document**

*Conventional approach 1:* Filter out other modalities besides unstructured text

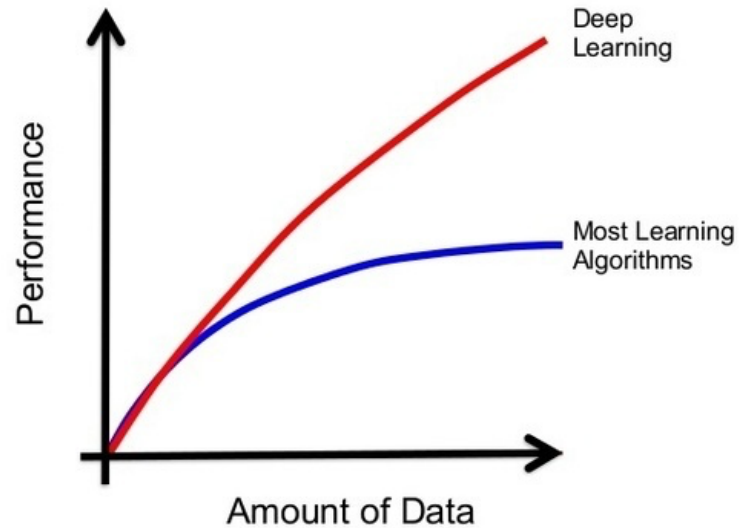*Conventional approach 2:* Limit the context scope to sentences or tables.

**Problem:** Misses important relations if you neglect multimodal information

Up to 97% missed relations!

# Deep learning is very successful in many domains



**Andrej Karpathy** [Follow]
Director of AI at Tesla. Previously Research Scientist at OpenAI and PhD student at Stanford. I like to train deep neural nets on large datasets.
Nov 11, 2017 · 8 min read

## Software 2.0

I sometimes see people refer to neural networks as just "another tool in your machine learn... work here or there, and som... ...ns. Unfortunately, Neural networ... of a fundamen...

## H2O Deep Learning beats MNIST

```
> install.packages("h2o")
> library(h2o)
> h2oServer <- h2o.init(ip="mr-0xd1", port=53322)
> train_hex <- h2o.importFile(h2oServer, "mnist/train.csv.gz")
> test_hex  <- h2o.importFile(h2oServer, "mnist/test.csv.gz")
> record_model <- h2o.deeplearning(x = 1:784, y = 785, data = train_h
                        activation = "RectifierWithDropo
                        epochs = 8000, l1 = 1e-5, input_
                        train_samples_per_iteration = -1
|                                               | 100%
> record_model@model$confusion
        Predicted
Actual   0    1    2    3    4    5    6    7    8    9 Error
  0    974    1    1    0    0    0    2    1    1    0 0.00612
  1      0 1135    0    1    0    0    0    0    0    0 0.00088
  2      0    0 1028    0    1    0    0    3    0    0 0.00388
  3      0    0    1 1003    0    0    0    3    2    1 0.00693
  4      0    0    0    0  971    0    0    0    6    0 0.01120
  5      2    0    0    5    0  882    1    1    1    0 0.01121
  6      2    3    0    1    1    2  949    0    0    0 0.00939
  7      1    2    6    0    0    0    0 1019    0    0 0.00875
  8      1    0    1    3    0    4    0    2  960    3 0.01437
  9      1    2    0    0    4    3    0    2    0  997 0.01189
```

**François Chollet** ✔
@fchollet [Following]

It is my impression that the world of deep learning *research* is starting to plateau. What's booming: deploying DL to real-world problems.

11:19 AM - 9 Sep 2017

186 Retweets 484 Likes

💬 26    🔁 186    ♡ 484

## Alibaba's artificial intelligence bot beats humans at reading in a first for machines

A deep neural network model developed by Alibaba has scored higher than humans in a reading comprehension test, paving the way for bots to replace people in customer service jobs

PUBLISHED : Monday, 15 January, 2018, 11:33am
UPDATED : Monday, 15 January, 2018, 12:17pm

### SQuAD
The Stanford Question Answering Dataset

## KEY MOMENTS IN DEEP-LEARNING HISTORY 2014-2016

**2014**
JANUARY
Google acquires DeepMind, a startup specializing in combining deep learning and reinforcement learning, for $600 million.

**2015**
DECEMBER
A team from Microsoft, using neural nets, outperforms a human on the ImageNet challenge.

**2016**
MARCH
DeepMind's AlphaGo, using deep learning, defeats world champion **Lee Sedol** in the Chinese game of go, four games to one.

LEE JIN-MAN—AP PHOTO

Can we take advantage of this powerful tool and apply it to our problem?

# Keys to utilizing deep learning



**How do we gather enough
labeled, richly formatted data?**

**How do we model the characteristics
of richly formatted data in DL?**

**Fonduer**
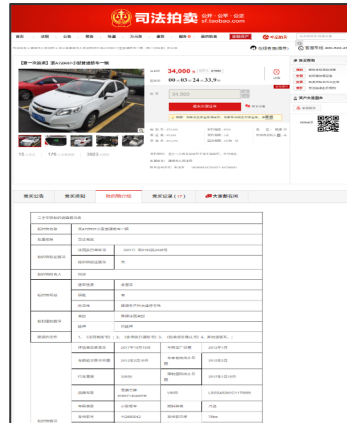A _weakly supervised_ deep learning framework for knowledge base construction from richly formatted data

# Fonduer in practice!



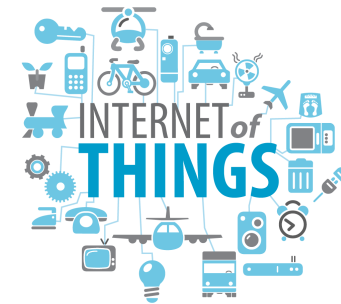Anti-Human Trafficking

Search Engine

Genome-wide
Association Studies

Internet
of Things

Paleontology

# Fonduer pipeline

# Generating richly formatted training data

# Multimodal weak supervision



**Weak supervision**: express any supervision signal via labeling functions to generate training data

```
# Check if current is in the same row with keyword `collector`
def in_the_same_row_with(candidate):
    if 'collector' in row_ngrams(candidate.current):
        return 1
    else: return -1
```
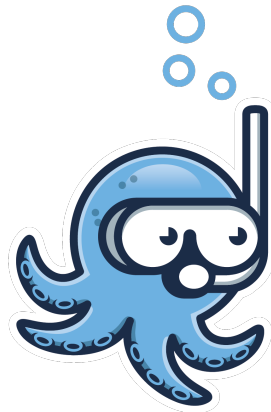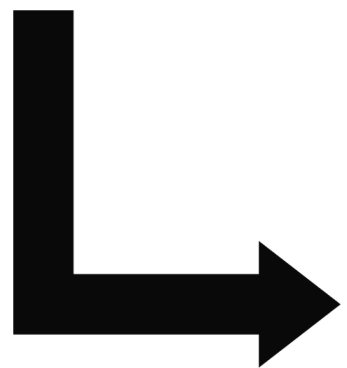
# Modeling Weak Supervision

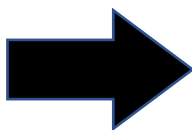| Doc. level Candidates | Multimodal Supervision | | |
|---|---|---|---|
| | Vertically aligned with 'Value' | Row ngrams contain 'mA' | 'current' in sentence |
| SMBT3904 *100* | ✗ | ∅ | ✓ |
| SMBT3904 *200* | ✓ | ✓ | ✗ |
| SMBT3904 *150* | ✓ | ✗ | ✗ |

∅ =Abstain

**Intuition**: Use agreements / disagreements to learn the accuracy of LFs without ground truth
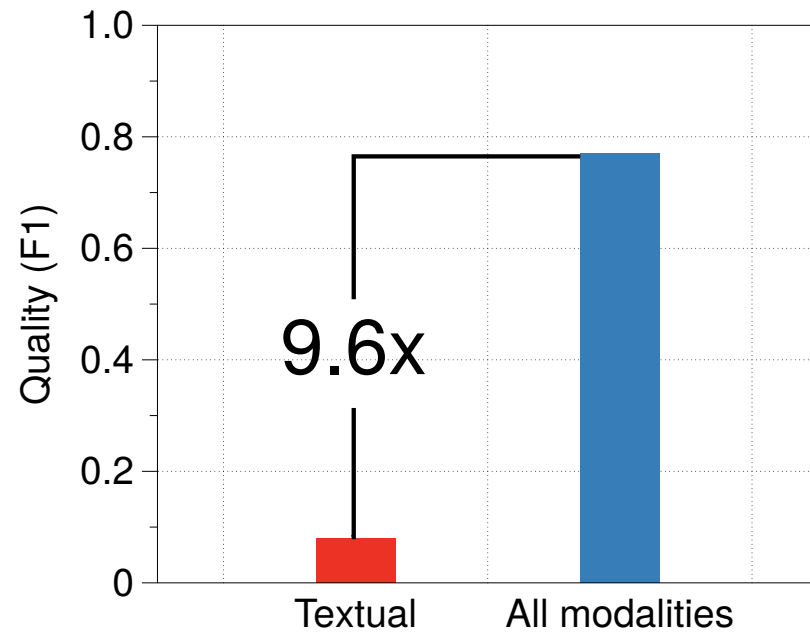
**Output**: Probabilistic Training Labels

| | | |
|---|---|---|
| SMBT3094 | *100* | **0.5** |
| SMBT3094 | *200* | **0.85** |
| SMBT3094 | *150* | **0.15** |

*Data programming/MeTal*

FONDUER

## For transistor datasheets…

Different supervision resources' effect



Modality distribution of supervision



**Users intuitively rely on multimodal information for supervision**

# Featurization and Classification for Richly Formatted Data

# LSTM for Textual Information

**Transistor Datasheet**
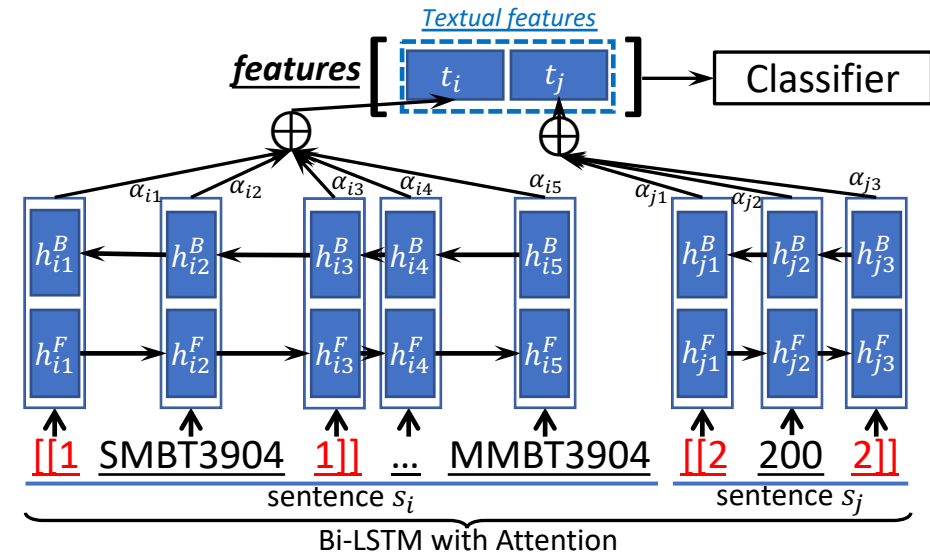
**SMBT3904**...MMBT3904

**NPN Silicon Switching Transistors**
- High DC current gain: 0.1 mA to 100 mA
- Low collector-emitter saturation voltage

**Maximum Ratings**

| Parameter | Symbol | Value | Unit |
|---|---|---|---|
| Collector-emitter voltage | $V_{CEO}$ | 40 | V |
| Collector-base voltage | $V_{CBO}$ | 60 | |
| Emitter-base voltage | $V_{EBO}$ | 6 | |
| Collector current | $I_C$ | 200 | mA |
| Total power dissipation | $P_{tot}$ | | mV |
| $T_S \leq 71°C$ | | 330 | |
| $T_S \leq 115°C$ | | 250 | |
| Junction temperature | $T_i$ | 150 | °C |
| Storage temperature | $T_{stg}$ | -65 ... 150 | |

LSTM excels at relation extraction from text
*Xu et al., 2015; Miwa et al., 2016; Zhang et al., 2016*



**Problem:** LSTM networks struggle to capture the multimodal characteristics of richly formatted data.

# Augmenting LSTM with Multimodal Features



**Transistor Datasheet**

**Font**: Arial; **Size:** 12; **Style**: Bold {**SMBT3904**...**MMBT3904**

**NPN Silicon Switching Transistors**
- High DC current gain: 0.1 mA to 100 mA
- Low collector-emitter saturation voltage

**Maximum Ratings**

| Parameter | Symbol | Value | Unit |
|---|---|---|---|
| Collector-emitter voltage | $V_{CEO}$ | 40 | V |
| Collector-base voltage | $V_{CBO}$ | 60 | |
| Emitter-base voltage | $V_{EBO}$ | 6 | |
| Collector current | | 200 | |
| Total power dissipation | $P_{tot}$ | | |
| $T_S \leq 71°C$ | | | |
| $T_S \leq 115°C$ | | 250 | |
| Junction temperature | $T_i$ | 150 | °C |
| Storage temperature | $T_{stg}$ | -65 ... 150 | |

*Same Font*

**Aligned**

**Font**: Arial; **Size:** 10

**Header: 'Value';
Row:** 2; **Column**: 3

We use the multimodal information stored in the *document* to extract basic multimodal features:
- ❑ **Structural features**
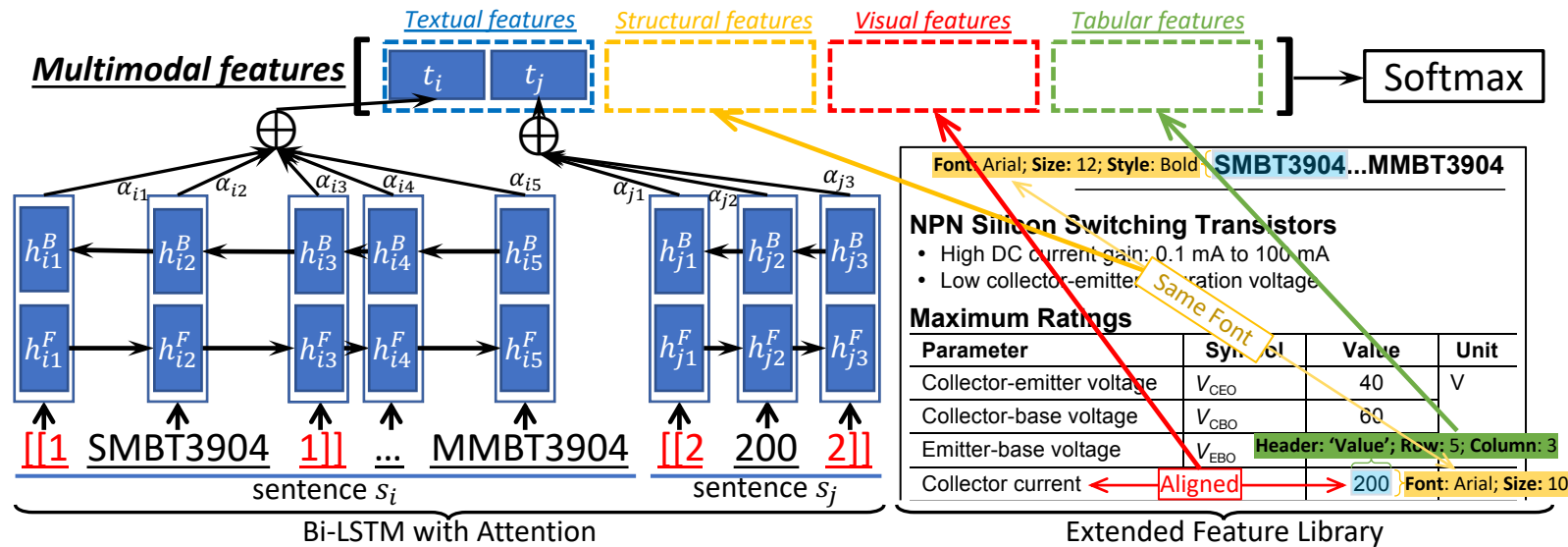- ❑ **Tabular features**
- ❑ **Visual features**

**Augmentation with multimodal features captures signals a traditional LSTM would miss.**
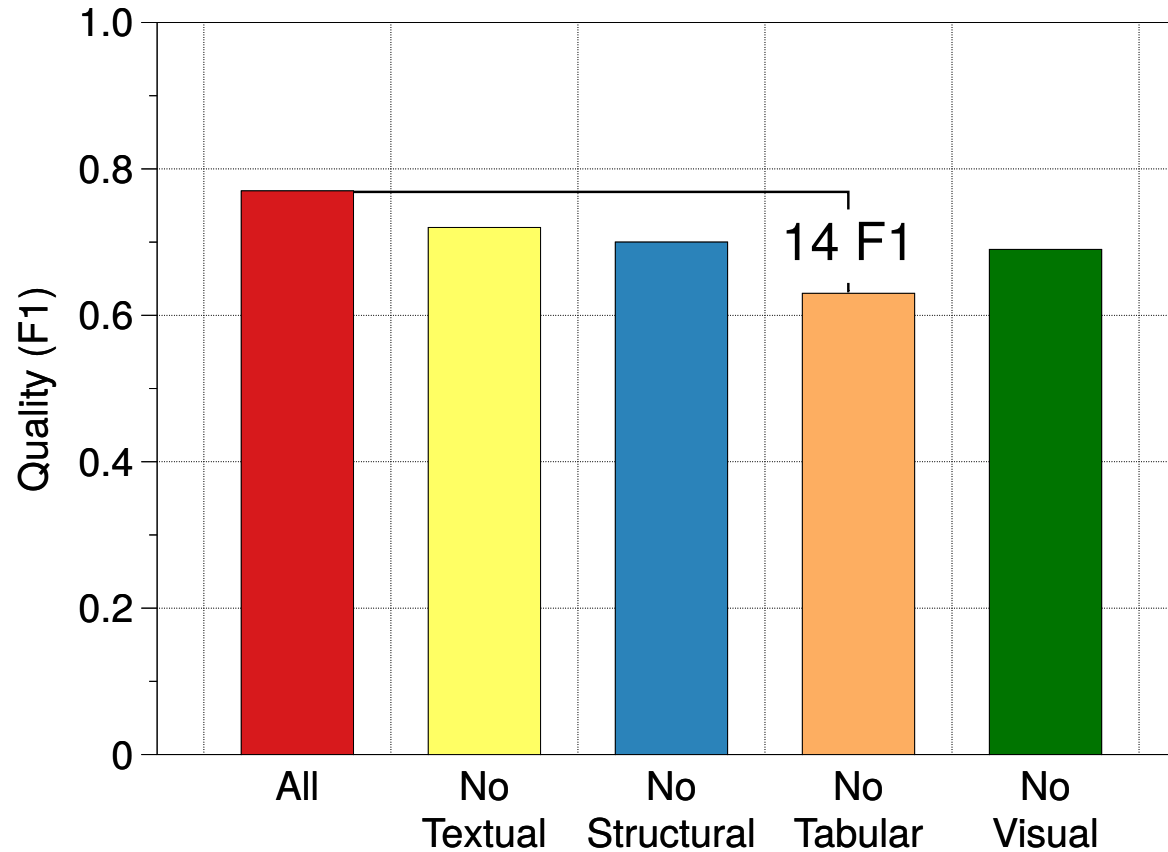
FONDUER

**Signals from different modalities can be useful to find the information.**



**Fonduer: a KBC system that takes advantage of both techniques to reason about all available signals.**

For transistor datasheets…



**Multimodal features significantly impact the quality of extraction**

# Fonduer in the wild
Empirical results & real-world uses

**FONDUER**

GWAS Catalog

**Fonduer**

Same set of documents

Human-created

10 years

1.0x extractions

Machine-created

**<6** months

**1.59x** extractions

Precision **0.89**
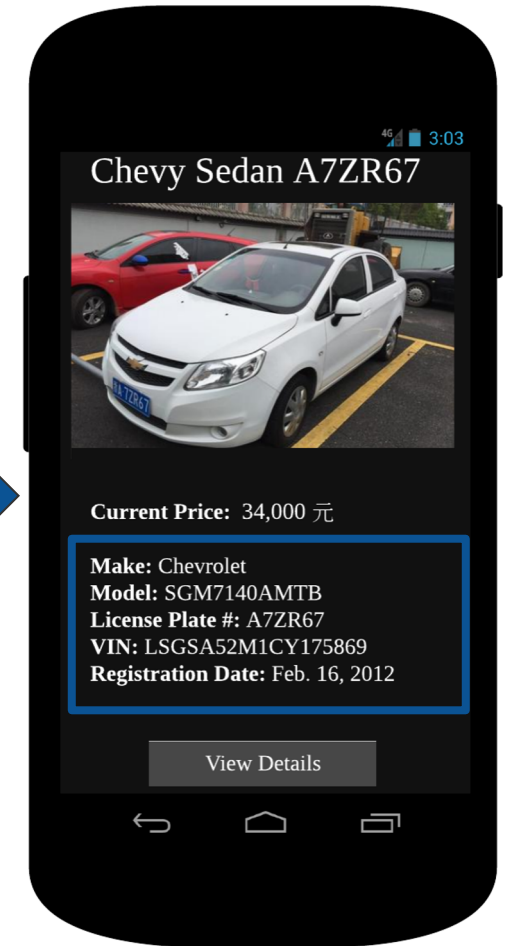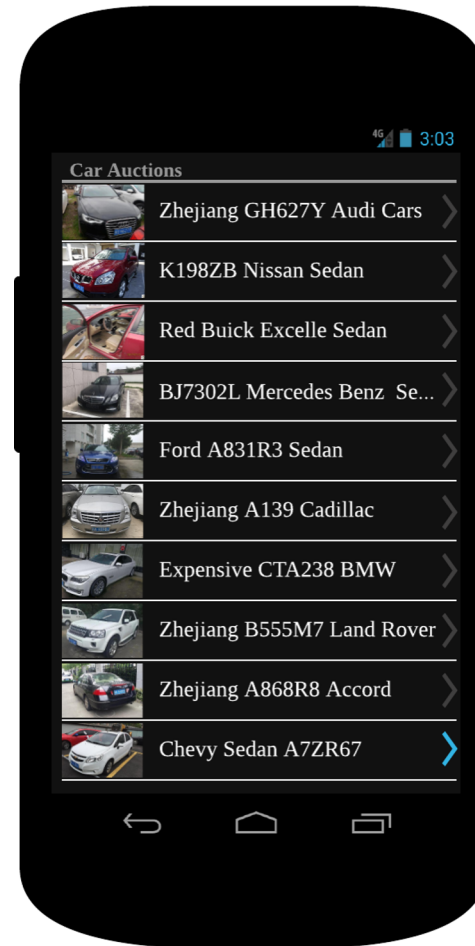
# How people use Fonduer in industry

**Input:** User-customized HTML auction pages → **Output:** Structured knowledge base

Extract key facts (make, model, license, etc.)

Improve auction search quality and UX



Fonduer

Alibaba.com

FONDUER

# Knowledge Base Construction from Richly Formatted Data

- Fonduer helps build high-quality KBC from richly formatted data

- Allows users to leverage multimodal signals

- Augments deep learning model with features from each data modality to achieve high quality

- Fonduer is supporting real world applications

**Thank you!**
**Sen Wu**
**(senwu@cs.stanford.edu)**

https://github.com/HazyResearch/fonduer