# The Joy of Character level NLP

—

Tal Perry - CEO LightTag
[tal@lighttag.io](mailto:tal@lighttag.io)
@thetalperry

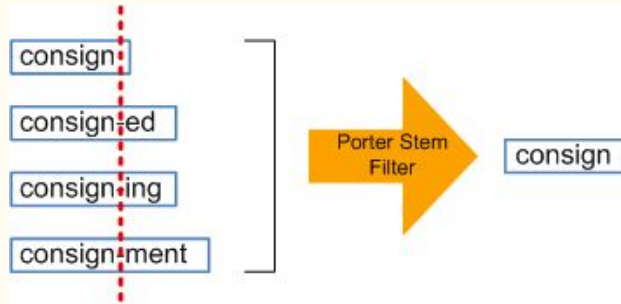# Character level NLP

All the cool kids are doing it

- Background
- Why ?
- Benefits
- Drawbacks
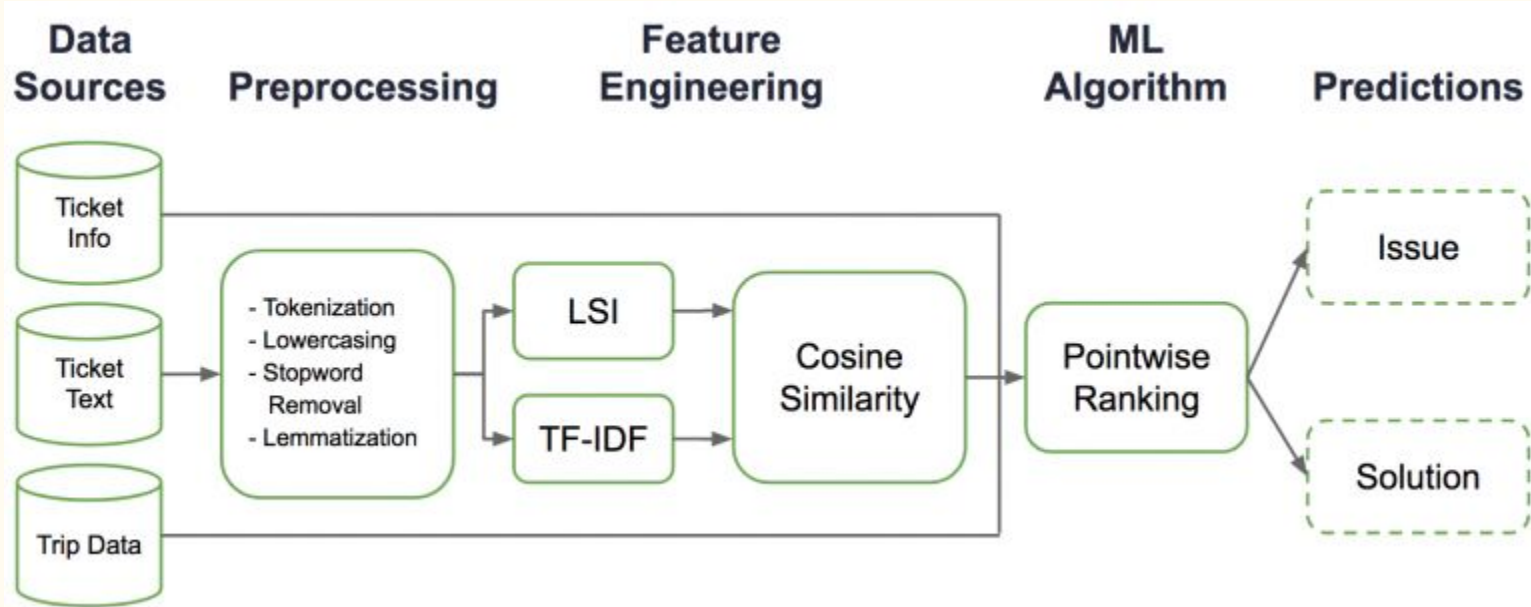- Example

# Background

# The Old *Old* Way

# The Old Old Way

# DL Insight 1: Word Vectors

0. *frog*
 1. frogs
 2. toad
 3. litoria
 4. leptodactylidae
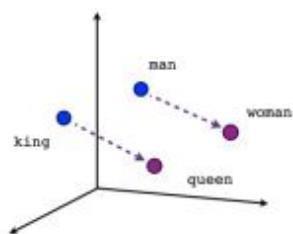 5. rana
 6. lizard
 7. eleutherodactylus

3. litoria

4. leptodactylidae

5. rana
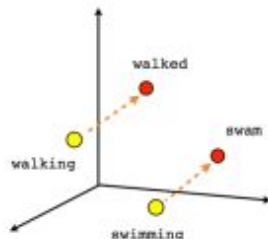
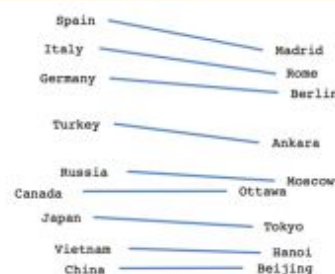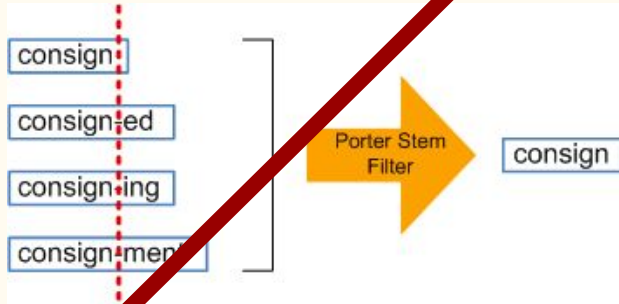7. eleutherodactylus

Male-Female

Verb tense

Country-Capital

# The Old *Old* Way

# Deep Learning made it easy to work with long sequences

```
1  import tensorflow as tf
2  lstm = tf.contrib.rnn.BasicLSTMCell(lstm_size)
```

# RNNs

```c
/*
 * Increment the size file of the new incorrect UI_FI
 * of the size generatively.
 */
static int indicate_policy(void)
{
  int error;
  if (fd == MARN_EPT) {
    /*
     * The kernel blank will coeld it to userspace.
     */
    if (ss->segment < mem_total)
      unblock_graph_and_set_blocked();
    else
      ret = 1;
    goto bail;
  }
  segaddr = in_SB(in.addr);
  selector = seg / 16;
  setup_works = true;
  for (i = 0; i < blocks; i++) {
    seq = buf[i++];
    bpf = bd->bd.next + i * search;
    if (fd) {
      current = blocked;
    }
  }
  rw->name = "Getjbbregs";
  bprm_self_clearl(&iv->version);
  regs->new = blocks[(BPF_STATS << info->historidac)]
  return segtable;
}
```

# The Old *Old* Way

# More RNNs

This is one of Crichton's best books. The characters of Karen Ross, Peter Elliot, Munro, and Amy are beautifully developed and their interactions are exciting, complex, and fast-paced throughout this impressive novel. And about 99.8 percent of that got lost in the film. Seriously, the screenplay AND the directing were horrendous and clea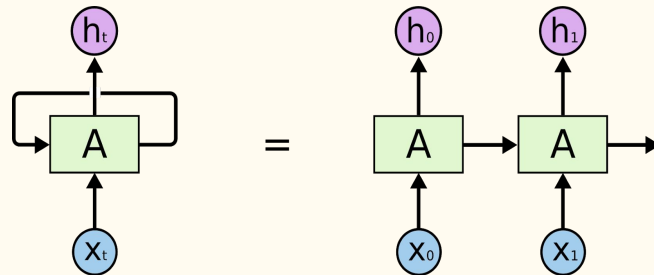rly done by people who could not fathom what was good about the novel. I can't fault the actors because frankly, they never had a chance to make this turkey live up to Crichton's original work. I know good novels, especially those with a science fiction edge, are hard to bring to the screen in a way that lives up to the original. But this may be the absolute worst disparity in quality between novel and screen adaptation ever. The book is really, really good. The movie is just dreadful.

# Why Characters ?

# But there are many "words"



Donald J. Trump ✔ @realDonaldT... · 49m ⌄
Despite the constant negative press covfefe

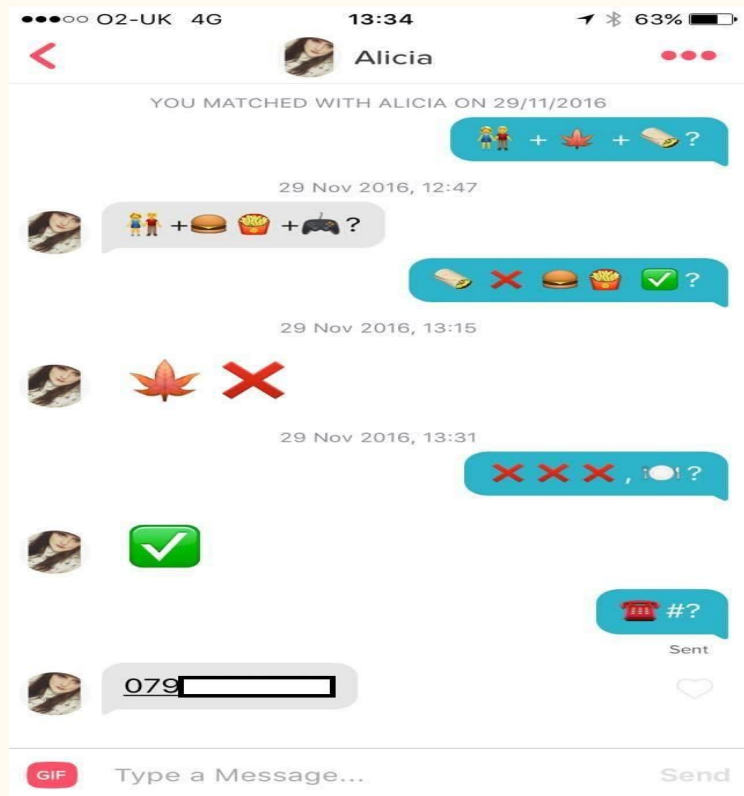↩ 15K    🔁 28.4K    ♥ 34.6K    ✉

# With many meanings

What has four letters, never has five letters but sometimes has nine letters.

# With many meanings

What has four letters, never has five letters but sometimes has nine letters?

# WTF ?

# Infinite

Vocabulary

- Expanding Vocabulary

- Many Professional Terms

- Lots of numbers / codes

- User generated data

# Infinite

Vocabulary

10/30/04: Cont. C&T. Drld cmt plug f/5693' to 5707'. Performed LOT to 960 psi = 15.6 ppg. Drld to 5875'.

# Infinite

## Vocabulary

# Infinite

Vocabulary

# Operational

## Simplicity

- No need to carry a dictionary

- Less preprocessing

- Less room for error

# In Action

**(a) Spelling mistakes**

| DE ori | Warum sollten wir nicht Freunde sei ? |
|---|---|
| DE src | Warum solltne wir nich Freunde sei ? |
| EN ref | Why should not we be friends ? |
| bpe2char | Why are we to be friends ? |
| char2char | Why should we not be friends ? |

**(b) Rare words**

| DE src | Siebentausendzweihundertvierundfünfzig . |
|---|---|
| EN ref | Seven thousand two hundred fifty four . |
| bpe2char | Fifty-five Decline of the Seventy . |
| char2char | Seven thousand hundred thousand fifties . |

**(c) Morphology**

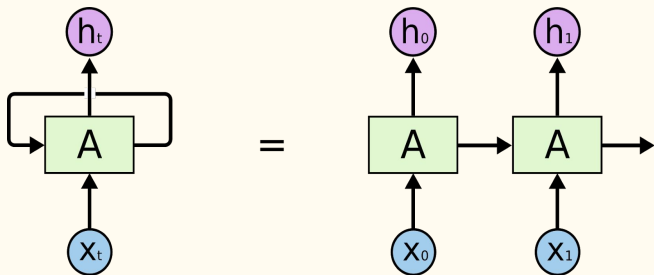| DE src | Die Zufahrtsstraßen wurden gesperrt , wodurch sich laut CNN lange Rückstaus bildeten . |
|---|---|
| EN ref | The access roads were blocked off , which , according to CNN , caused long tailbacks . |
| bpe2char | The access roads were locked , which , according to CNN , was long back . |
| char2char | The access roads were blocked , which looked long backwards , according to CNN . |

**(d) Nonce words**

| DE src | Der Test ist nun über , aber ich habe keine gute Note . Es ist wie eine Verschlimmbesserung . |
|---|---|
| EN ref | The test is now over , but i don't have any good grade . it is like a worsened improvement . |
| bpe2char | The test is now over , but i do not have a good note . |
| char2char | The test is now , but i have no good note , it is like a worsening improvement . |

# Drawbacks

# Long Sequences

Mean more computation



- RNNs are O(n) in sequence length

- The average english word has ~5 letters
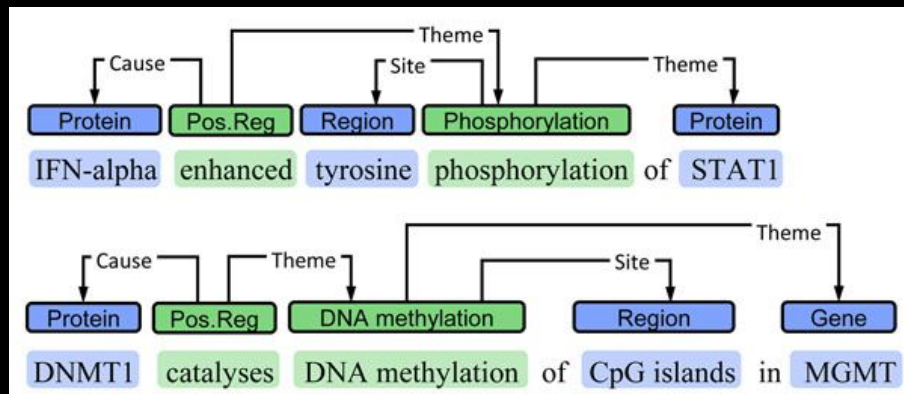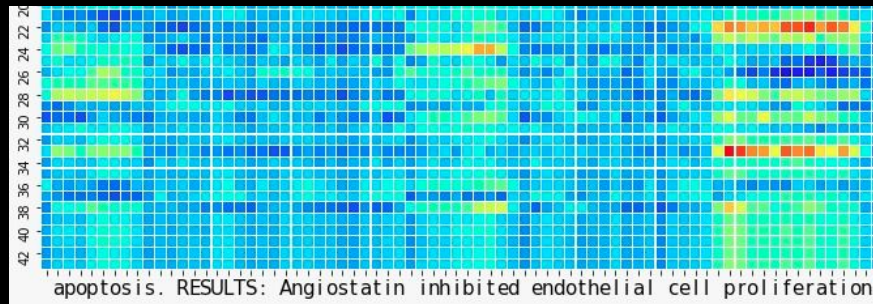
- So 5x more expensive

# Meaningless

## Semantic Unit

- Characters have no semantic meaning

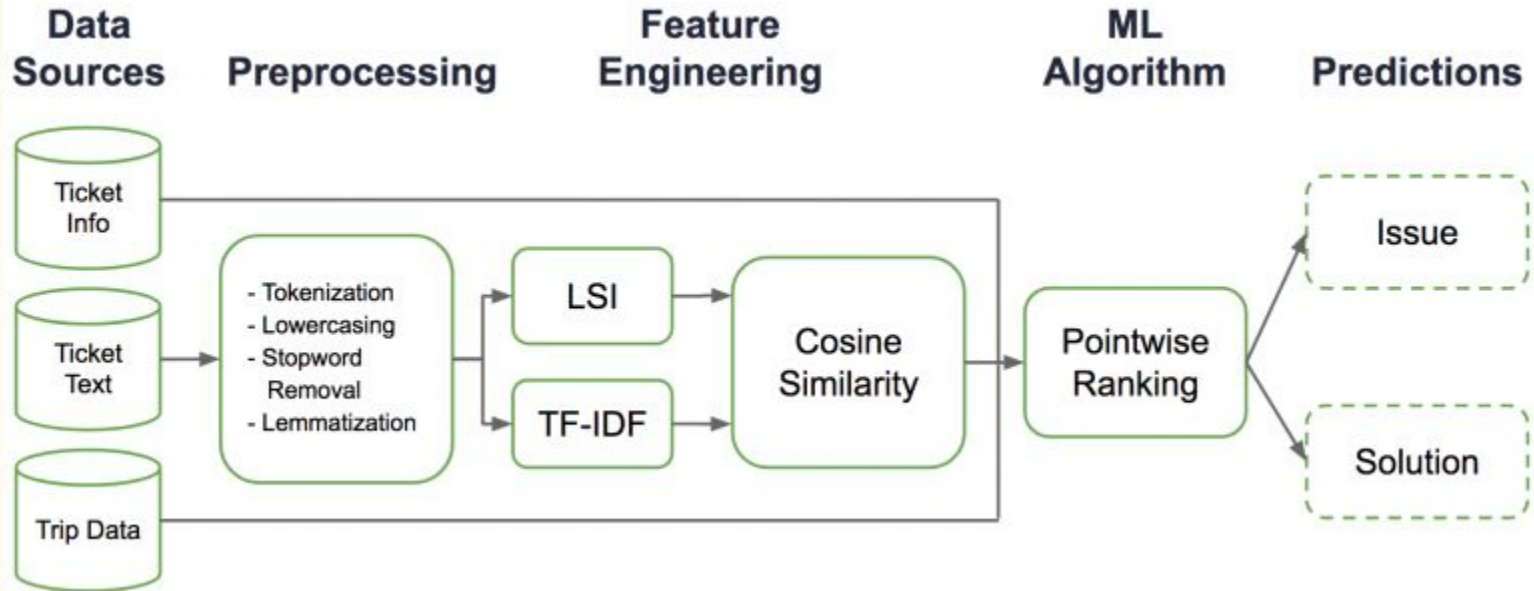- Its harder to form relations between a bunch of meaningless things

# Meaningless

## Semantic Unit

# Example

*https://github.com/LightTag/BibSample*

# Conclusion

# The End

# Questions ?