



Machine learned search: setting up a production pipeline

Zalando SE

Maximilian Werk
Senior Research Engineer

20-01-2020







Machine learned search: setting up a production pipeline

Zalando SE

Maximilian Werk
Senior Research Engineer

20-01-2020



Programming vs. Software Engineering

Who is a Software Engineer?

Who does programming in their day-to-day work?

Who has a machine learning background?

Bekleidung

Schuhe

Sport

Accessoires

Premium

Beauty

Sale

Geschenkgutscheine

'boho kleid'

Sortieren  Größe  Marke  Farbe  Preis  Obermaterial  Muster  Mehr Filter

33 Artikel

Gesponsert



WEEKEND MaxMara 374,95 €
TORNADO - Strickkleid - anthrazit



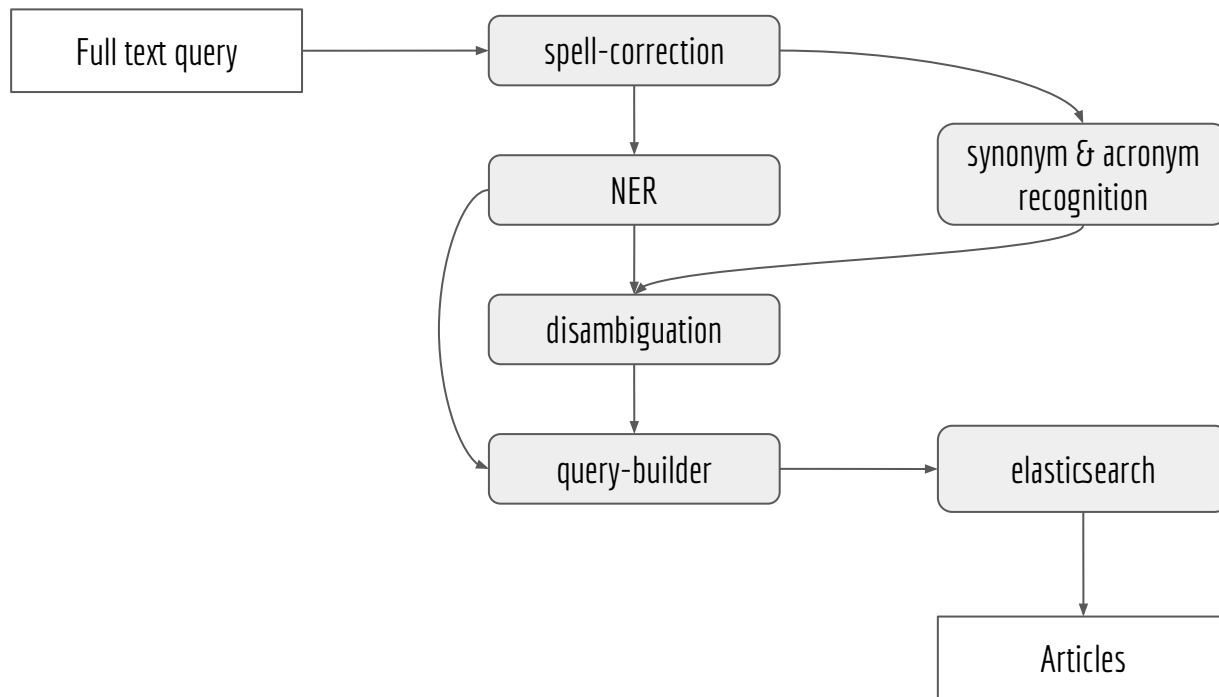
YAS 49,99 €
Freizeitkleid - sudan brown



-10%

Nly by Nelly 54,95 €
BOHO TUNIC DRESS - Freizeitkleid... 49,49 €

Our Information Retrieval Pipeline



Failing classical information retrieval

“pullover patchwork”



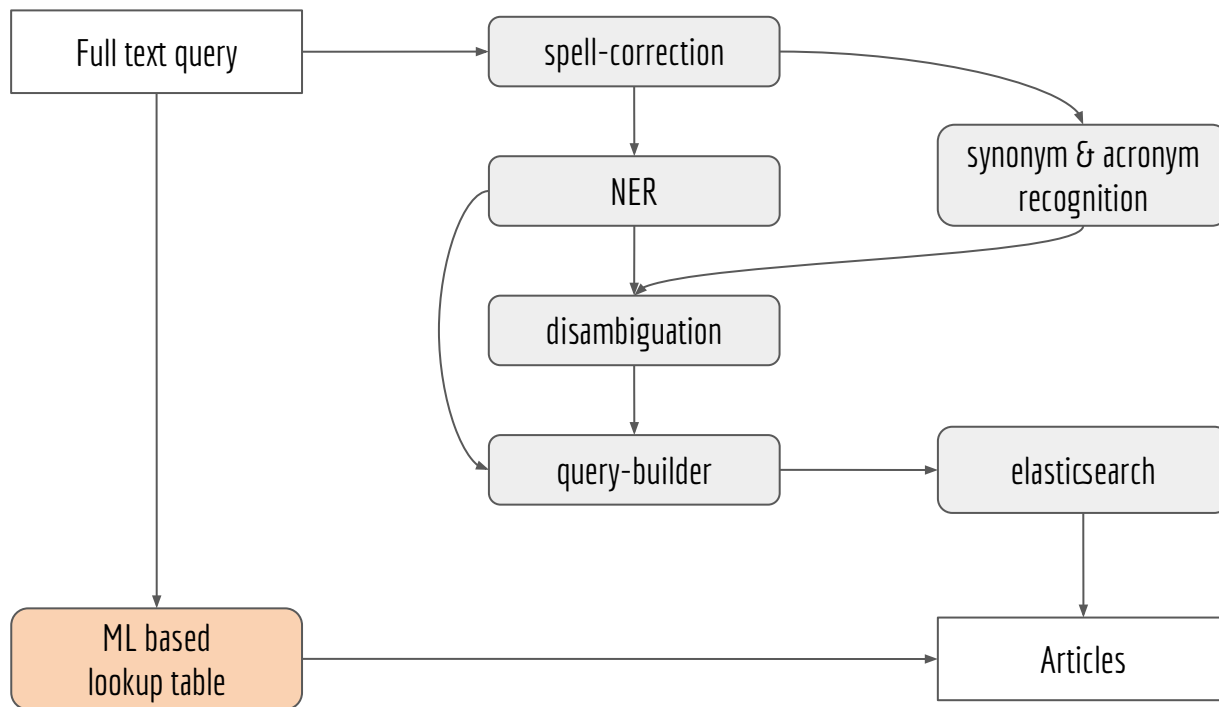
“top figurumspielend”



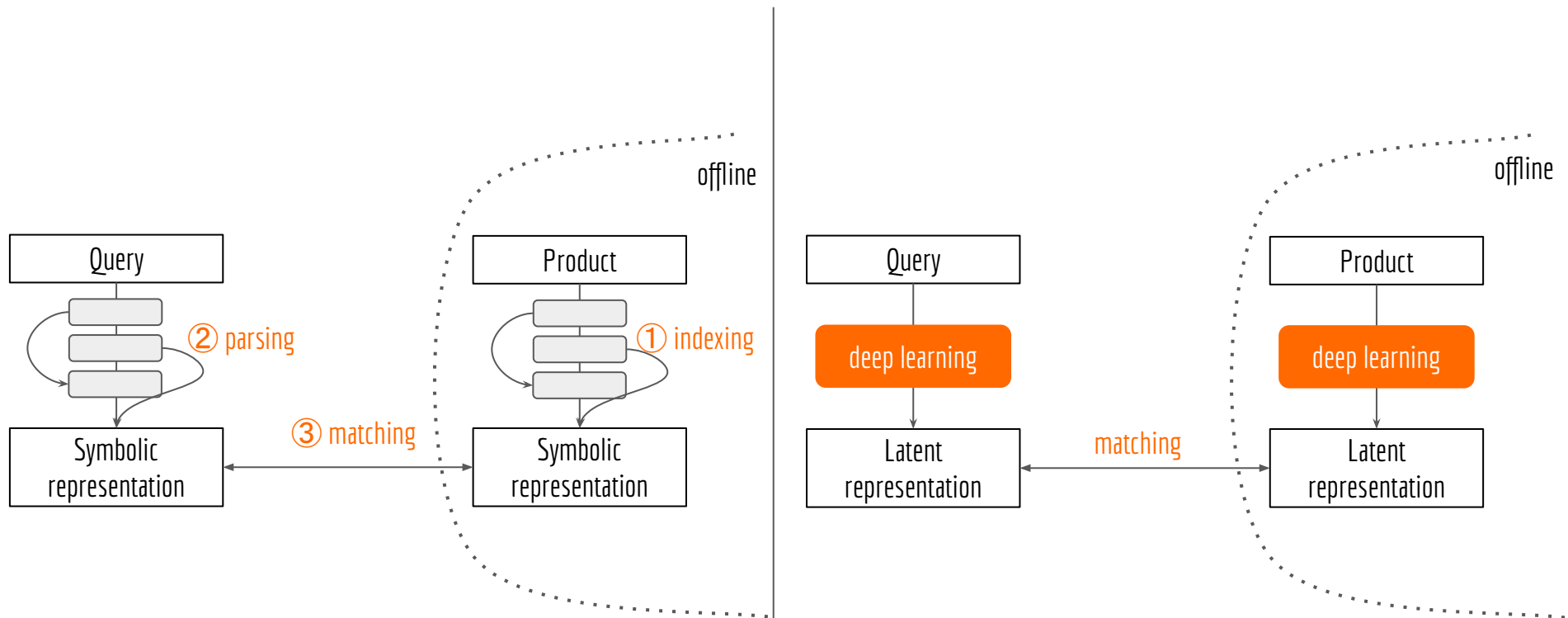
“abendkleid tattoospitze”



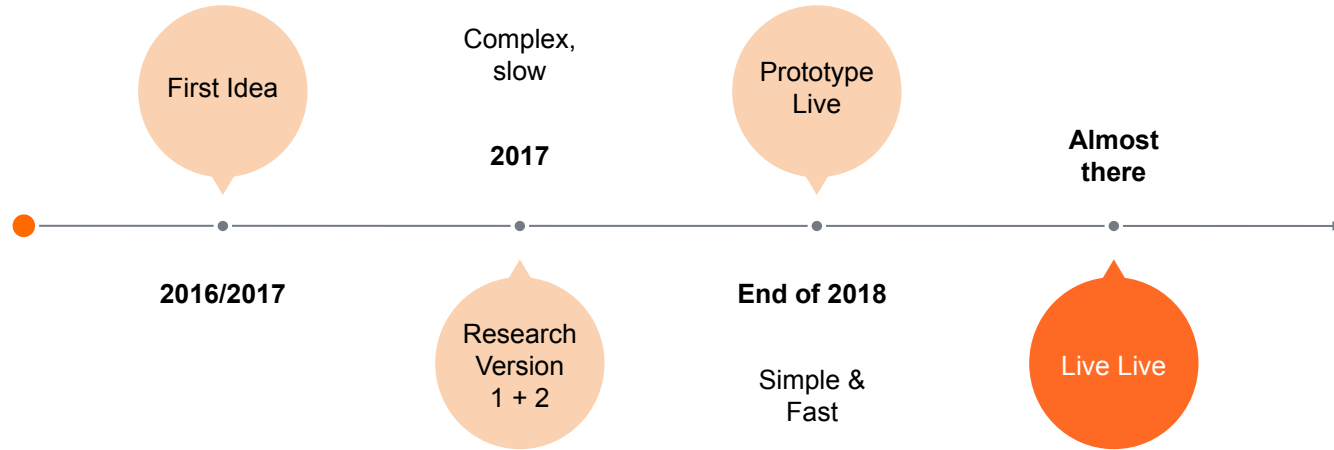
Adding ML based solution



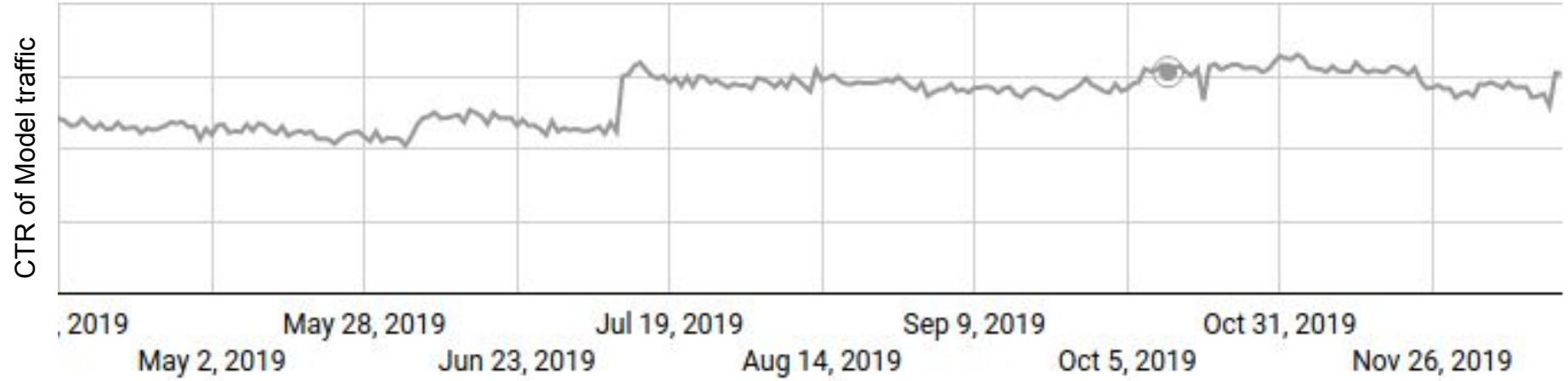
Classical system vs. end-to-end product search system



History



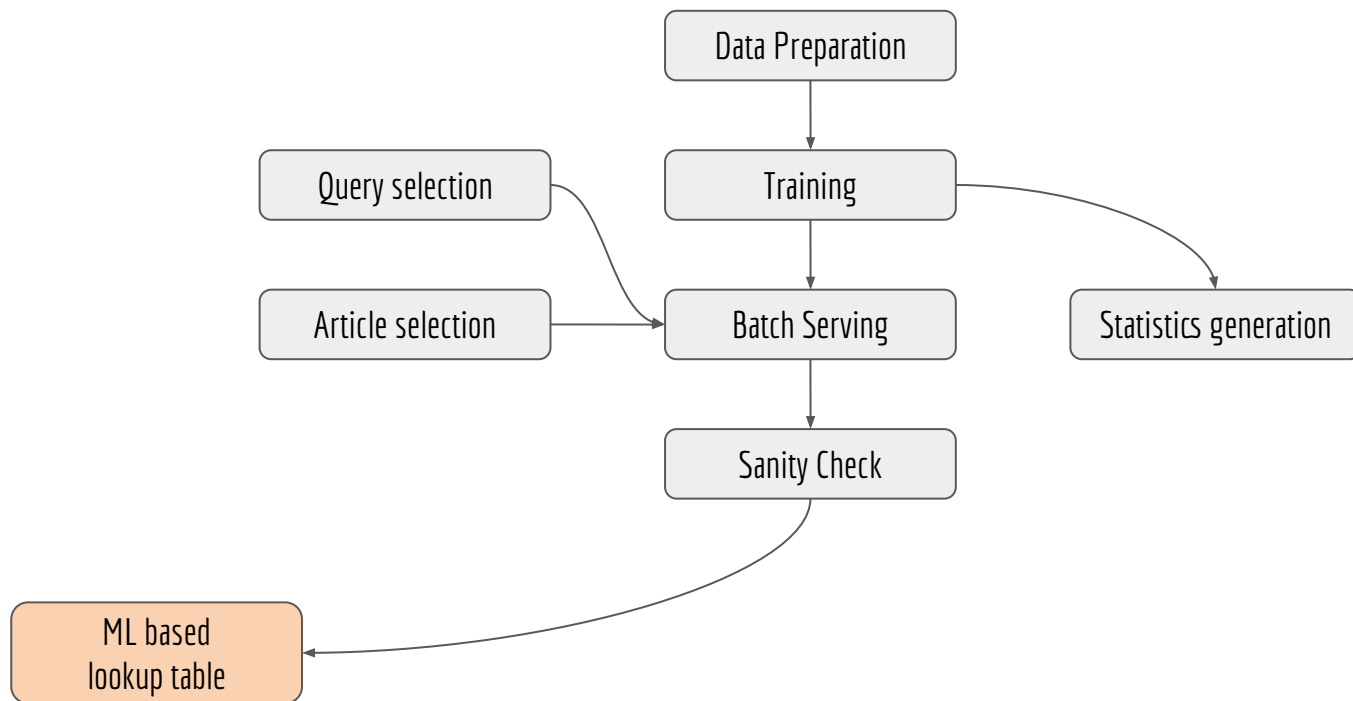
Model degradation



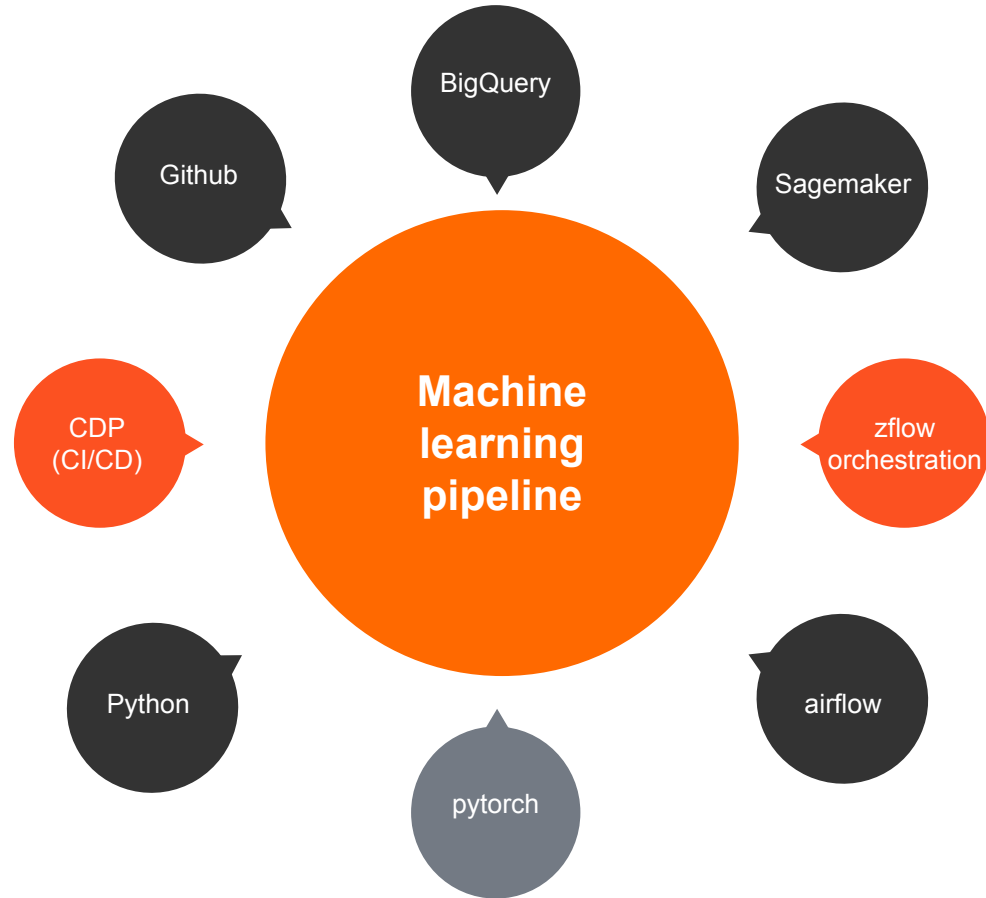
Building a pipeline



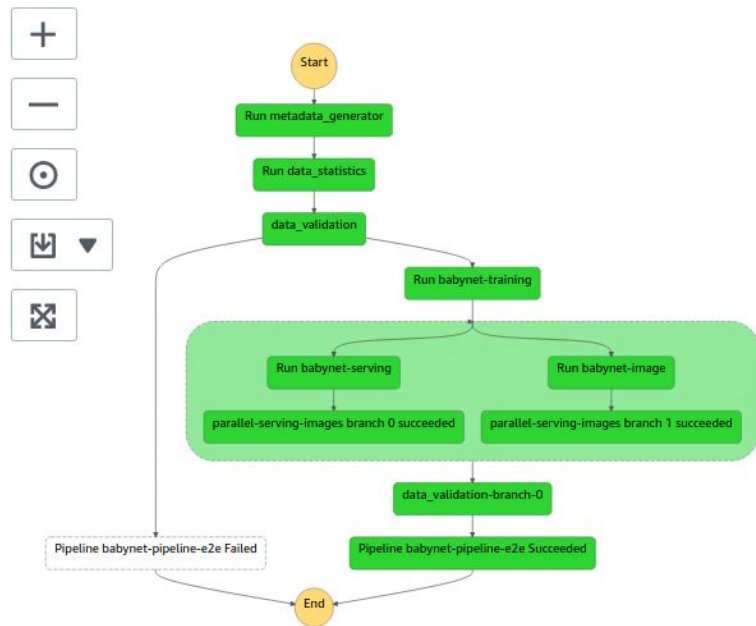
ML pipeline



Our Technology Stack



Visual workflow



■ In Progress ■ Succeeded ■ Failed ■ Cancelled ■ Caught Error

Code

Step details

Name Type

Run babynet-image Task

Status

✓ Succeeded

Resource

[REDACTED]

► Input

► Output

► Exception

One-time job vs. continuous development

“Obvious decisions”

Complex model vs. simple model

Batch vs. live

Manual vs. automated

Training cadence (daily, weekly, monthly, irregularly)

One-time job vs. continuous development

“Hidden decisions”

Fast now vs. fast, maintainable, robust later

Scripting vs. Software Development

Talk, talk talk!



Testing is hard



Configuration is complex



Follow standards & best practices



Do clean code

```
def p(i):  
    s = i.split("|")  
    d = s[0].lower()  
    # Cleaning the price string into python number representation and converting it to cents  
    p = int(float(s[1].replace(",", ".").replace("€", "")) * 100)  
    return d, p
```

BÄÄÄH

```
from collections import namedtuple  
  
Article = namedtuple("Article", ["description", "price"])  
  
def parse_article(article_string):  
    splitted = article_string.split("|")  
    description = splitted[0].lower()  
  
    cleaned_price_string = splitted[1].replace(",", ".").replace("€", "")  
    price_in_cents = int(float(cleaned_price_string) * 100)  
  
    return Article(description, price_in_cents)
```

Nice

About good (ML) code

Correct

Simple functions

Written for others to read

Accessible business logic

Pipeline steps are independently executable

Producing good code

- 1) Feature
- 2) Correct
- 3) Readable
- 4) Simple
- 5) Readable
- 6) Feature is still correct?
 - a) No => go to 1)
 - b) Yes => Happy days!

Be a scout: Leave the code cleaner than you found it.

Future work

Improve Model

Monitoring

Model influences training data

Add more use cases

Replace tradition IR search

Takeaways

Simple model does the job

Pipeline building takes a lot of time

Train code craftsmanship



Maximilian Werk



Search - Team Lens
Senior Research Engineer

maximilian@zalando.de

Twitter: @maintainable_ds



We hire:

[Senior Search Engineer](#)

[Principal Research Engineer - Search](#)

[Principal Product Manager - Search](#)

20-01-2020

