# About us

- Spinout of UCL's Computer Science department, specialising in computational advertising and electric commerce

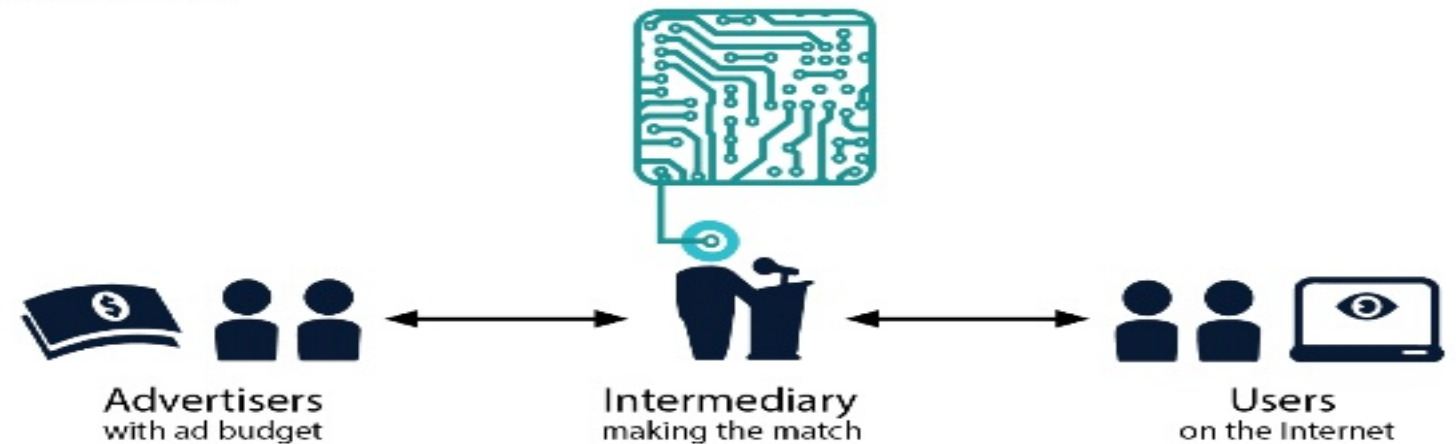| | Platform development begins | | First long term contract (Telefonica) | | | 11 clients | |
|---|---|---|---|---|---|---|---|
| **JAN'15** | **MAR'15** | **FEB'16** | **AUG'16** | **FEB'17** | **AUG'17** | | **FEB'18** |
| UCL spinout created | | First PoC revenue | | First DSP client roll out | | | Pre-A |

- Proved our technology in the ad tech industry w/clients such as Beeswax & Telefonica
- Currently process over 3TB per day, containing tens of billions of daily user events, across tens of millions mobile profiles, spanning 5 countries.
- We work with DSPs/SSPs/exchanges & telcos w/over 85% accuracy & less than 10ms latency
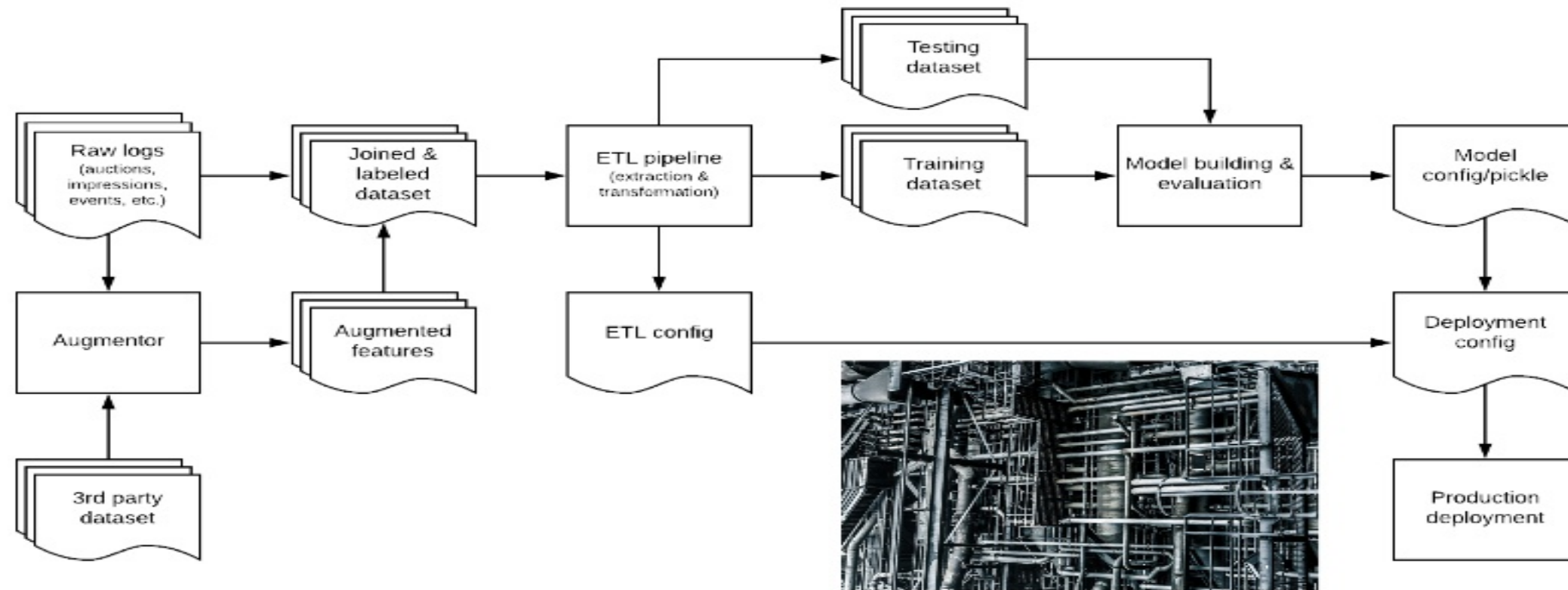
# What do we do

- FRAUD
    - Is the user human?
    - Up to 40% of ads are not shown to humans

- ACTIONS
    - How likely is the user to click on the ad or install an app or register?

- PRICING
    - How much should pay for this impression?

- RELEVANCE
    - Is this user my target audience?
    - How do I find more of the same users?

# Real-time Decision Making

- Real-Time
  - Thousands of QPS
  - 99.9% response under 10ms

- Bidding
  - User response prediction (e.g., CTR prediction)
  - Bid price

- Optimisation
  - ROI
  - Volume (i.e., budget spent)

**Advertisers**
with ad budget

**Intermediary**
making the match

**Users**
on the Internet

# An end-to-end pipeline

# Feature Engineering

# Feature Engineering contd.

```
1 [
2      "timestamp$month$7",
3      "timestamp$day$1",
4      "timestamp$weekday$4",
5      "timestamp$hour$0",
6      "timestamp$minute$0",
7      "exchange$nexage",
8      "bidrequest$app$publisher$ext$nex_data_rights$0",
9      "bidrequest$app$publisher$id$16797",
10     "bidrequest$app$publisher$name$24/7 apps",
11     "bidrequest$app$domain$247apps.com",
12     "bidrequest$app$name$24/7 apps-playtube free-android",
13     "bidrequest$app$bundle$com.tfsapps.playtube2",
14     "bidrequest$app$cat$iab19-17",
15     "bidrequest$app$cat$iab1-5",
16     "bidrequest$app$ext$nex_sdkv$5.3.0-c3980670.a",
17     "bidrequest$app$id$55290",
18     "bidrequest$app$storeurl$https://play.google.com/store/apps/details?id=com.tfsapps.pla
19     "bidrequest$regs$coppa$0",
20     "bidrequest$imp$pmp$deals$id$1426189778844608480",
21     "bidrequest$imp$bidfloor$1.0",
22     "bidrequest$imp$ext$nex_screen$0",
23     "bidrequest$imp$instl$0",
24     "bidrequest$imp$banner$h$50",
25     "bidrequest$imp$banner$pos$1",
26     "bidrequest$imp$banner$battr$3",
27     "bidrequest$imp$banner$battr$4",
28     "bidrequest$imp$banner$battr$5",
29     "bidrequest$imp$banner$battr$8",
30     "bidrequest$imp$banner$battr$9",
31     "bidrequest$imp$banner$battr$12",
32     "bidrequest$imp$banner$api$5",
33     "bidrequest$imp$banner$w$320",
34     "bidrequest$imp$banner$btype$1",
35     "bidrequest$at$2",
36     "bidrequest$device$language$en",
37     "bidrequest$device$make$samsung",
38     "bidrequest$device$lmt$1",
```
```
1 [
2             42239,
3             83074,
4            140934,
5            208266,
6            244091,
7            244443,
8            305412,
9            328341,
10           352227,
11           414817,
12           424476,
13           438697,
14           512487,
15           512867,
16           598740,
17           604956,
18           608432,
19           675206,
20           706406,
```

# Challenge 1

- How to deal with arbitrary fields in unstructured logs?



Expansion to year/month/day/hour etc. required

Augmentation opportunities

Deeply nested

Multi-items in value

Some fields should be dropped

# Challenge 2

- How to guarantee the feature extraction/augmentation consistency?

```
>>> pp.pprint(user_agent_parser.Parse('Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)
Chrome/69.0.3497.100 Safari/537.36'))
{   'device': {   'brand': None, 'family': 'Other', 'model': None},
    'os': {   'family': u'Windows',
              'major': u'10',
              'minor': None,
              'patch': None,
              'patch_minor': None},
    'string': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/69.0.3497.100 Saf
ari/537.36',
    'user_agent': {   'family': 'Chrome',
                      'major': '69',
                      'minor': '0',
                      'patch': '3497'}}
```

It'll be a huge headache
if happens on important features

```
In [2]: user_agent_parser.Parse('Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome
   ...: /69.0.3497.100 Safari/537.36')
Out[2]:
{'device': {'brand': None, 'family': 'Other', 'model': None},
 'os': {'family': 'Windows 10',
 'major': None,
 'minor': None,
 'patch': None,
 'patch_minor': None},
 'string': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/69.0.3497.100 Safari
/537.36',
 'user_agent': {'family': 'Chrome',
 'major': '69',
 'minor': '0',
 'patch': '3497'}}
```

# Challenge 3

- How to make the ETL process portable?

# Challenge 4

- How to do it fast enough?
  - Hundreds of thousands of QPS
  - 10-15ms round trip time
  - Overhead for API & decoding (e.g., protobuf)
  - Cost?
  - It's common to implement the prediction functions in a different language (than python)

# OpenETL

- Tree traversal
  - A recursive function
  - Deals with both structured and unstructured input requests

- Libs + Configuration
  - Build libs for multiple programming language
  - Load configurations at runtime
  - Different levels of tests to guarantee consistency

- Micro services architecture; containerize:
  - I/O
  - Common ETL
  - Experiment control
  - Specific transformation & model & stacking

# Alternatives

- **Featuretools**
  - A framework to perform automated feature engineering. It excels at transforming temporal and relational datasets into feature matrices for machine learning.
  - Featuretools is intended to be run on datasets that can fit in memory on one machine.

- TransmogrifAI
  - An end-to-end AutoML library for structured data written in Scala that runs on top of Apache Spark. It was developed with a focus on accelerating machine learning developer productivity through machine learning automation, and an API that enforces compile-time type-safety, modularity, and reuse.

- bonobo
  - A lightweight Extract-Transform-Load (ETL) framework for Python 3.5+

- https://www.featuretools.com
- https://transmogrif.ai
- https://www.bonobo-project.org

# Extraction

- Operators
  - Object traverse
    - Lists & dicts
    - Optional depth limit
  - Split
  - Exclude
  - Augment
    - Internal & external datasource
  - Evaluate
    - Essentially `eval()`
    - E.g., converting timestamps

# Augmentation

- Examples
  - doc2vec for a given corpus
  - Historical CTR/CVR
  - First-party user data (e.g., abandoned shopping cart value)
  - Time + location -> weather

- Integration
  - As dictionary
  - Real-time API

# Transformation

- Operators
  - CountVectorizer
  - HashingVectorizer
  - Bucketizer
  - MinMaxScaler
  - PolynomialFeatures



- If necessary, trained in Apache Spark
  - For many transformation `fit()` is expensive but `transform()` is cheap
  - E.g., `OpenETLCountVectorizer.copy_from_spark()`

- Rosicrucian Digest on Alchemy, https://www.rosicrucian.org/rosicrucian-digest-alchemy

# Optimisation

- Higher-level APIs to manipulate the ETL pipeline steps
  - Step selection in training -> step importance
  - Optional priority field
    (dropping steps/features when performance degrades)

- Cython for python, later other programming languages
  - Golang
  - Java



- WikiMedia Commons

# A Real-world Example

- Input

  – Customer defined

  – Text based

  – JSON format

  – Requires further processing

# A Real-world Example, contd.

Feature extraction
by traversing
JSON/pyobj tree

```json
1    [
2        "created_at": "2018-09-01 15:09:19",
3        "steps": [
4            {
5                "namespace": "extract",
6                "class": "ExtractPythonObject",
7                "arguments": {
8                    "delimeter": "$"
9                }
10           },
11           {
12               "namespace": "extract",
13               "class": "SplitFeature",
14               "arguments": {
15                   "seperator": "|",
16                   "feature": "user$categories",
17                   "delimeter": "$"
18               }
19           },
20           {
21               "namespace": "extract",
22               "class": "SplitFeature",
23               "arguments": {
24                   "seperator": "|",
25                   "feature": "video$title",
26                   "delimeter": "$"
27               }
28           },
29           {
30               "namespace": "extract",
31               "class": "SplitFeature",
32               "arguments": {
33                   "seperator": "|",
34                   "feature": "video$type",
35                   "delimeter": "$"
36               }
37           },
```
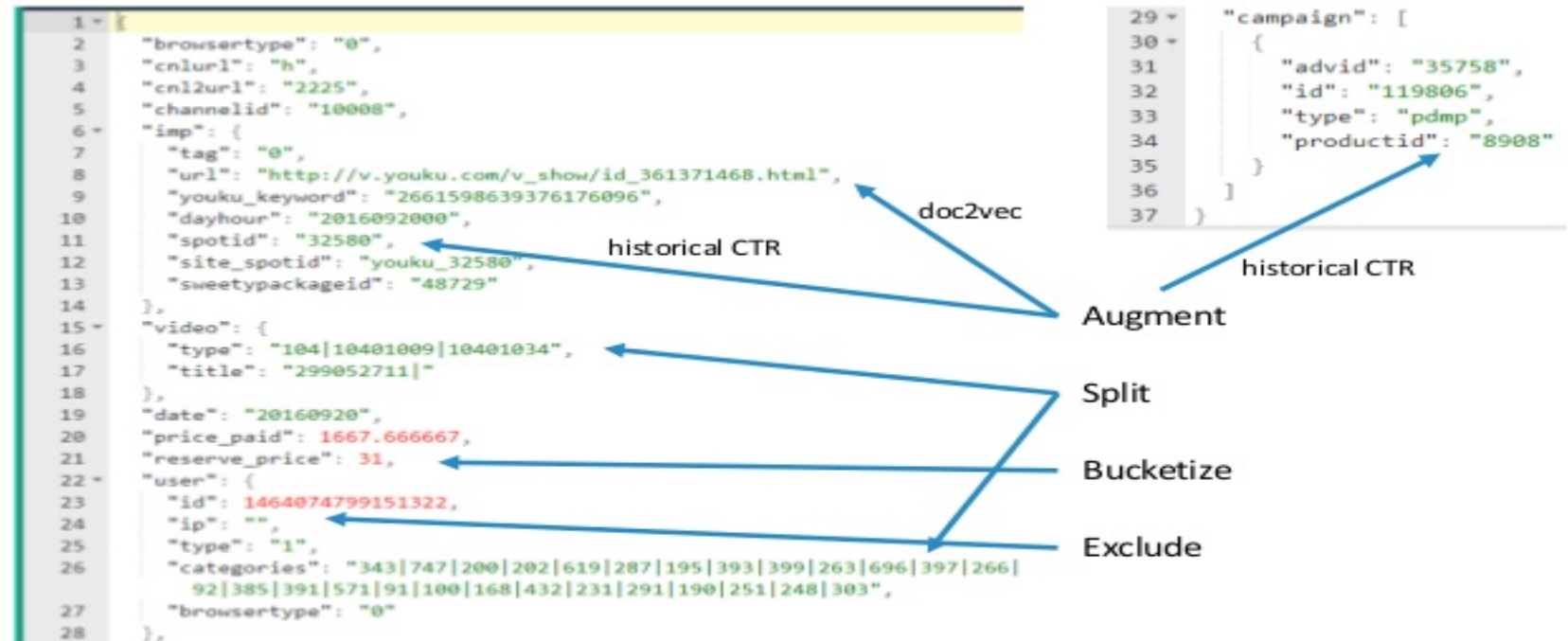
```json
38           {
39               "namespace": "extract",
40               "class": "ExcludeFeature",
41               "arguments": {
42                   "feature": "user$id"
43               }
44           },
45           {
46               "namespace": "extract",
47               "class": "ExcludeFeature",
48               "arguments": {
49                   "feature": "user$ip"
50               }
51           },
```

```json
52           {
53               "namespace": "extract",
54               "class": "AugmentFeature",
55               "arguments": {
56                   "feature": "url",
57                   "vocabulary": [
58                       {
59                           "http://www.abc.com": [0,1,2,3,4,5,"..."]
60                       },
61                       "..."
62                   ],
63                   "default_value": [0,0,0,0,"..."]
64               }
65           },
66           {
67               "namespace": "transform",
68               "class": "CountVectorizer",
69               "arguments": {
70                   "vocabulary": [
71                       "..."
72                   ],
73                   "size": 197805,
74                   "binary": true
75               }
76           }
77       ],
78       "name": "Demo ETL model"
79   }
```

Embedded dictionary
for augmentation

Vectorisation
by OneHotEncoding

# A Real-world Example, contd.

- Output:
  - Dense / sparse vector: size, indices, values
  - JSON/CSV/Parquet
  - Optional "label" field
- Utilities for format conversion
  - org.apache.spark.ml.linalg.SparseVector
  - scipy.sparse.csr_matrix
  - tf.SparseTensor
  - etc.

```
 1   {
 2      "size": 197805,
 3      "indices": [
 4         0,
 5         1,
 6         2,
 7         3,
 8         4,
 9         8,
10        14,
11        15,
12        "...",
13        22316
14      ],
15      "values": [
16         1,
17         1,
18         1,
19         1,
20         1,
21        "...",
22         1
23      ]
24   }
```

# Thank you!

- Questions?

- We are hiring!

- Shuai Yuan, VP Data Science, MediaGamma
- shuai.yuan@mediagamma.com