



OpenETL for Real-time Decision Making

Shuai Yuan, VP Data Science, MediaGamma

#SAISEnt7

About us

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



- 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?

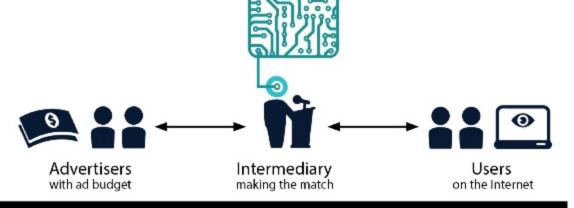
RELEVANCE

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



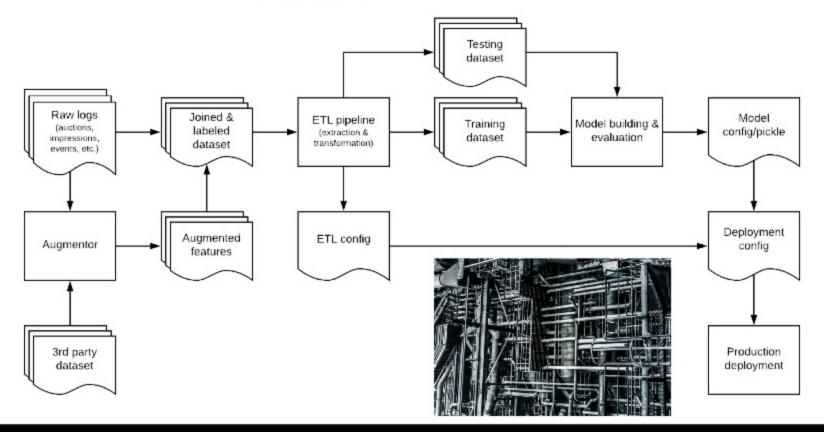
Real-time Decision Making

- 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)





An end-to-end pipeline





Feature Engineering

```
"timestamp": 1467331224000,
              "exchange": "Nexage", 
"bidRequest": {
                          "app": {
                                      "publisher": {
                                                "ext": {"nex_data_rights": 0}, "id": "19982",
                                                 "name": "myYearbook.com"
                                      "domain": "meetme.com",
                                     "name": "nyyearbook Android",
"bundle": "com.myyearbook.m",
"cat": ["IABIA"],
                                     "ext": {"nex_sdkv": "6.1.8-5323db4.a"},
"id": "70578",
"storeurl": "http://wex.neetne.com/"
                       },
"regs": {"coppa": 0),
"imp": [{
    "pep": {
    "deals": [
18
19
21
                                                           {"id": "1462976892784335237"},
                                                            ("id": "1426863662579673123"),
                                                            ("1d": "1435781799382281797")
                                                           ("1d": "1439493878945656671")
28
29
30
                                     ),
"bidfloor": 2.4,
                                     "displaymanagerver": "6.1.8-5323db4.a", 
"displaymanager": "millennial",
31
                                     "ext": {"nex_screen": 0},
                                    "instl": {
"banner": {
    "h": 50,
                                     "instl": 0,
33
                                                "battr": [1, 2, 3, 4, 8, 9, 10],
                                                "api": [5],
                                                "w": 320,
```

```
"timestamp$month$7",
           "timestamp$day$1",
           "timestamp$weekday$4",
           "timestamp$hour$0",
           "timestamp$minute$0",
           "exchange$nexage",
           "bidrequest$app$publisher$ext$nex_data_rights$0",
           "bidrequest$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$0"
           "bidrequest$imp$pmp$deals$id$1426189778844608480",
20
21
           "bidrequest$imp$bidfloor$1.0"
           "bidrequest$imp$ext$nex_screen$0",
22
23
           "bidrequest$imp$inst1$0"
24
           "bidrequest$imp$banner$h$50"
25
           "bidrequest$imp$banner$pos$1",
           "bidrequest$imp$banner$battr$3"
26
27
           "bidrequest$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$samsung",
           "bidrequest$device$lmt$1",
```



Feature Engineering contd.

```
"timestamp$month$7",
           "timestamp$day$1",
           "timestamp$weekday$4",
           "timestamp$hour$0",
           "timestamp$minute$0",
           "exchange$nexage",
           "bidrequest$app$publisher$ext$nex_data_rights$0",
           "bidrequest$app$publisher$id$16797"
           "bidrequest$app$publisher$name$24/7 apps",
           "bidrequest$app$domain$247apps.com",
           "bidrequest$app$name$24/7 apps-playtube free-android",
           "bidrequest$app$bundle$com.tfsapps.playtube2",
           "bidrequest$app$cat$iab19-17",
           "bidrequest$app$cat$iab1-5",
           "bidrequest$app$ext$nex_sdkv$5.3.0-c3980670.a",
           "bidrequest$app$id$55290",
           "bidrequest$app$storeurl$https://play.google.com/store/apps/details?id=com.tfsapps.pla
           "bidrequest$regs$coppa$0",
           "bidrequest$imp$pmp$deals$id$1426189778844608480",
           "bidrequest$imp$bidfloor$1.0",
           "bidrequest$imp$ext$nex_screen$0",
           "bidrequest$imp$inst1$0",
           "bidrequest$imp$banner$h$50".
           "bidrequest$imp$banner$pos$1",
           "bidrequest$imp$banner$battr$3"
           "bidrequest$imp$banner$battr$4",
           "bidrequest$imp$banner$battr$5",
           "bidrequest$imp$banner$battr$8",
           "bidrequest$imp$banner$battr$9"
           "bidrequest$imp$banner$battr$12",
           "bidrequest$imp$banner$api$5",
33
           "bidrequest$imp$banner$w$320",
           "bidrequest$imp$banner$btype$1",
35
           "bidrequest$at$2",
           "bidrequest$device$language$en",
37
           "bidrequest$device$make$samsung",
           "bidrequest$device$lmt$1".
```

```
42239,
            83074,
            140934.
            208266,
           244091,
           244443,
            305412,
           328341,
10
           352227,
11
           414817,
12
           424476.
13
           438697,
           512487,
15
           512867,
16
           598740,
17
            604956,
18
            608432.
19
           675206,
           706406,
```



Challenge 1

How to deal with arbitrary fields in unstructured logs?

```
Expansion to year/month/day/hour etc. required
"timestamp": 1467331224000,
"exchange": "Nexage"
"bidRequest": (
       "app": {
               "publisher": {
                                                                               Augmentation opportunities
                      "ext": {"nex_data_rights": 0}, 
"id": "19982",
                       "name": "myVearbook.com"
                                                                               Deeply nested
              "name": "myyearbook Android"
              "bundle": "com.myyearbook
              "cat": ["IAB14"]
                                                                               Multi-items in value
                                                                                                                                                                          "browsertype": "8",
              "id": "79578",
                                                                                                                                                                          "cnlurl": "h".
              "storeurl": "http://www.meetme.com/"
                                                                                                                                                                          "tag": "0",
      "regs": {"copps": 0},
"imp": [{
                                                                                                                                                                        "url": "http://v.youku.com/v_show/id_361371468.html",
                                                                               Some fields should be dropped
                                                                                                                                                                         "tanx_crowd": null,
                                                                                                                                                                          "baidu_usercategory"
                                                                                                                                                                            "343 747 200 202 619 287 195 393 399 263 696 397 266 92 385 391 571 91
                                                                                                                                                                            |100|168|432|231|291|190|251|248|303",
                              ("id": "1462976892784335237"),
                                                                                                                                                                          "advid": "35758",
                              ("id": "1426863662579673123"),
                                                                                                                                                                          "youku_keyword": "2661598639376176896",
                              ("id": "1435781799382281797"),
                              ("1d": "1439493878945656671")
                                                                                                                                                                         "dayhour": "2816892888",
                                                                                                                                                                         "spotid": "32588"
                                                                                                                                                                         "cn12ur1": "2225"
               "bidfloor": 2.4,
                                                                                                                                                                          "sweetypackageid": "48729",
              "displaymanagerver": "6.1.8-5323db4.a", 
"displaymanager": "millennial",
                                                                                                                                                                          "video_type": "104|10401009|10401034",
                                                                                                                                                                          "date": "20160920".
              "ext": {"nex_screen": 0}, 
"instl": 0,
                                                                                                                                                                          "price_paid": 1667.666667,
                                                                                                                                                                          "productid": "8988",
                                                                                                                                                                          "campaigntype": "pdmp"
                                                                                                                                                                          "visitorid": "1464874799151322",
                                                                                                                                                                          "campaignid": "119806"
                      "battr": [1, 2, 3, 4, 8, 9, 10],
                                                                                                                                                                          "site_spotid": "youku_32580",
                                                                                                                                                                         "channelid": "10005",
                      "btype": [1]
                                                                                                                                                                          "usertype": "1",
                                                                                                                                                                          "reserve_price": 31,
              "1d": "fb2d45c6-1655-6e6c-865a-e710ac7608e3-1"
                                                                                                                                                                          "video_title": "299052711|"
```

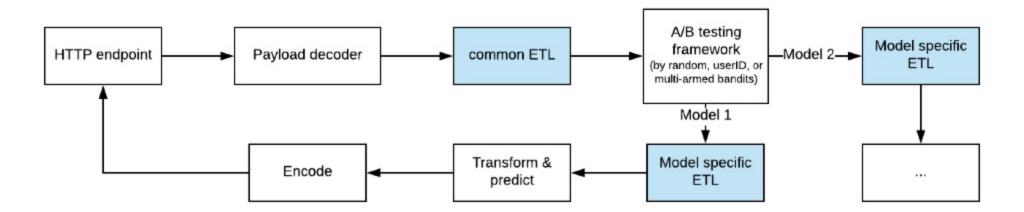
Challenge 2

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},
             'family': u'Windows',
                                                                                                                         It'll be a huge headache
             'major': u'10', <
             'minor': None,
                                                                                                                         if happens on important features
             'patch': None,
             'patch_minor': None},
   'string': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/60.3497.100 Sat
ari/537.36'.
   'user_agent': {
                     'family': 'Chrome',
                                              2 user_agent_parser.Parse(
                      'major': '69',
                      'minor': '0',
                      'patch': '3497'}
                                                                               Other', 'model': None),
                                             device': {'brand': None, 'family':
                                            '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.

TransmogrifAl

 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.

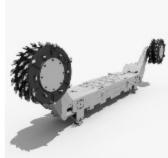
- A lightweight Extract-Transform-Load (ETL) framework for Python 3.5+
- https://www.featuretools.com
- https://transmogrif.ai
- https://www.bonobo-project.org



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

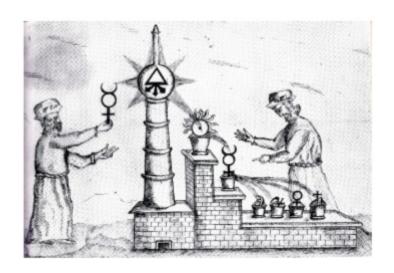




Transformation

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



A Real-world Example

Input

- Customer defined
- Text based
- JSON format
- Requires further processing

```
29 +
                                                                                         "campaign": [
                                                                                  30 *
      "browsertype": "0",
      "cnlur1": "h",
                                                                                 31
                                                                                             "advid": "35758",
      "cn12ur1": "2225".
                                                                                             "id": "119806",
                                                                                 32
      "channelid": "10008",
                                                                                  33
                                                                                             "type": "pdmp",
      "imp": {
                                                                                             "productid": "8908"
                                                                                 34
        "tag": "0",
                                                                                 35
        "url": "http://v.youku.com/v_show/id_361371468.html",
                                                                                  36
        "youku keyword": "2661598639376176096",
                                                                  doc2vec
       "dayhour": "2016092000",
                                                                                 37
11
        "spotid": "32580",
                                         historical CTR
12
        "site_spotid": "youku_32580"
                                                                                          historical CTR
13
        "sweetypackageid": "48729"
14
                                                                              Augment
15 -
      "video": {
        "type": "104|10401009|10401034",
16
        "title": "299052711|"
17
                                                                              Split
18
19
      "date": "20160920",
      "price_paid": 1667.666667,
20
21
      "reserve_price": 31,
                                                                              Bucketize
22 *
23
        "id": 1464074799151322,
        "ip": "",
24
                                                                              Exclude
25
        "type": "1",
        "categories": "343|747|200|202|619|287|195|393|399|263|696|397|266|
          92 385 391 571 91 100 168 432 231 291 190 251 248 303",
        "browsertype": "0"
27
28
```



A Real-world Example, contd.

Feature extraction by traversing JSON/pyobj tree

```
39
      "created_at": "2018-09-01 15:09:19",
      "steps": [
                                              40
                                              41 -
          "namespace": "extract",
                                              42
          "class": "ExtractPythonObject",
                                              43
          "arguments": {
            "delimeter": "$"
                                              45 +
                                              46
10
                                              47
11 -
                                              48 -
12
          "namespace": "extract".
                                              49
13
          "class": "SplitFeature",
                                              50
14 *
          "arguments": {
            "seperator": "|",
15
16
            "feature": "user$categories",
17
            "delimeter": "$"
18
19
20 =
21
          "namespace": "extract",
          "class": "SplitFeature",
22
23 -
          "arguments": {
            "seperator": "|"
24
            "feature": "video$title".
25
            "delimeter": "$"
27
28
29 *
30
          "namespace": "extract",
31
          "class": "SplitFeature",
32 -
          "arguments": (
            "seperator": "|"
33
            "feature": "video$type",
34
35
            "delimeter": "$"
36
```

```
"namespace": "extract",
  "class": "ExcludeFeature",
  "arguments": {
    "feature": "user$id"
}
},
{
    "namespace": "extract",
    "class": "ExcludeFeature",
    "arguments": {
        "feature": "user$ip"
}
}
```

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



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.

```
2
       "size": 197805.
       "indices": [
 3 +
         14.
10
11
         15,
12
13
        22316
14
15 +
       "values": [
16
17
        1,
18
19
21
22
23
24
```

Thank you!

- Questions?
- We are hiring!
- Shuai Yuan, VP Data Science, MediaGamma
- shuai.yuan@mediagamma.com

