# AWS Ground Truth
☐ *is*
☐ *is not*
# A silver bullet

What AWS has done for you
and what you must be prepared to do

# PAVEL DERBAN

*Lead BI Business Analyst*
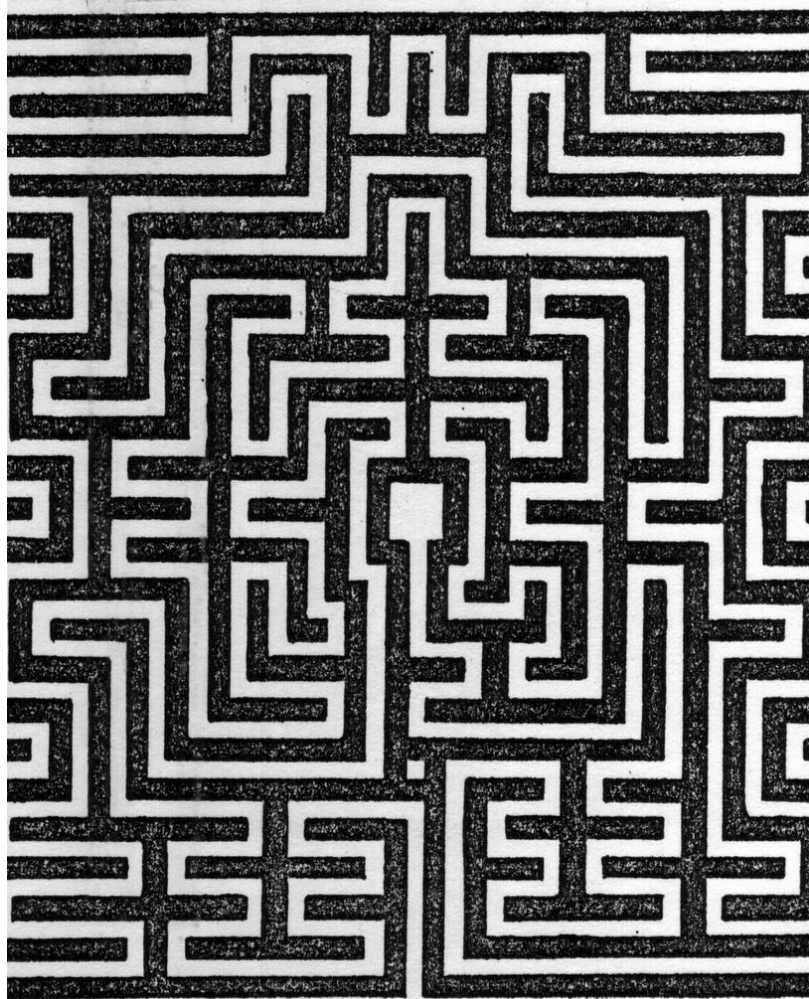
Fascinated by ML

# AGENDA

# TERMINOLOGY INTRO

- **Supervised Machine Learning** is a learning from labeled training data

- **Training Data Set** is labeled data used to train machine learning algorithms

- **Named Entity (NE)** is a proper name referring to important business object (company name, person name, security name)

- **Coreference/Anaphora** is an expression referencing to the mentioned named entity (he, the suspect, the president)

- **Taxonomy** is treelike structure of classes of objects within a domain

- **Knowledge Graph** is a brain-like, structured database that stores facts with flexible, bi-directional relationships

- **Knowledge Graph triple feed (RDF triple)** is a set of three entities that codifies a statement about semantic data in the form of subject–predicate–object expressions

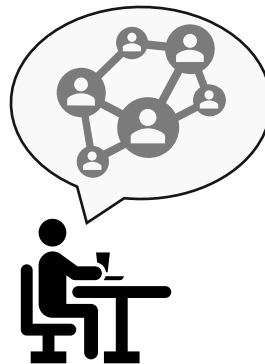# BUSINESS CASE

- EPAM client's business:         news aggregation,
                                              information extraction and provision
- Final product:                           access to knowledge graph containing information
- Client's customer:                   investment banks
- Business Value:                        investment risk assessment
- Enabler product:                      person-risk relation automated discovery
- Enabler of enabler product:    training data set

*Where is ML?*

*Where is Ground Truth?*

*Where is Cloud?*

# ML TASK FORMULATION

- Business Value:
  automated news classification,
  human effort reduction

- Input:
  news feed, xml files
- Preprocessed (3ʳᵈ Party):
  named entity recognition (NER),
  co-reference resolution (CR)

- Processing:
  taxonomy class (Risk) prediction,
  NE-Risk relation prediction

- Output:
  graph-ready triplet (quadruplet):
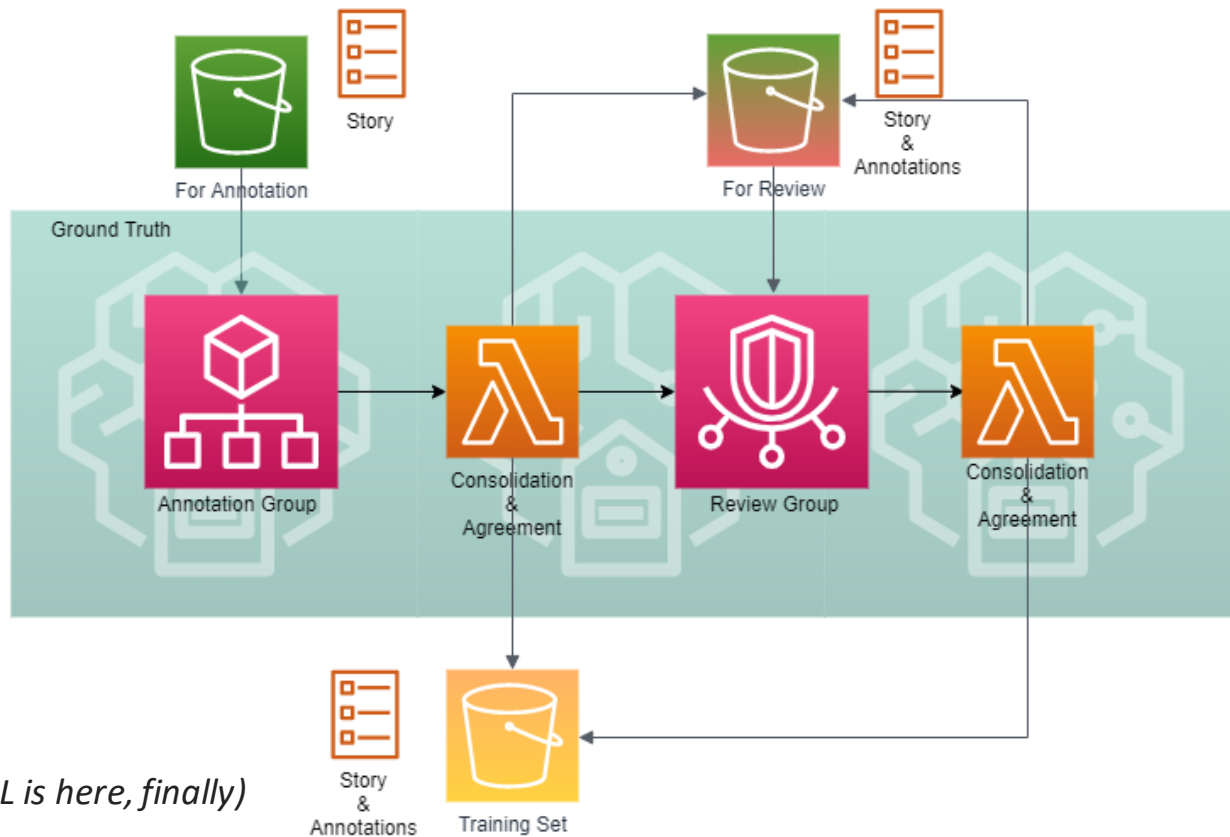  "subject-relation-risk-(source)"

```
News  >>>  NER + CR  >>>  Risk        >>>  Relation    >>>
                         Prediction        Prediction
```

# TRAINING SET TASK

- **Business Value:**          provide human labeled "ground truth" to train model, reduce cost of training set through provision of convenient annotation tool

- **Input:**          batch of xml files, custom schema
- **NE/Coreference resolved**          NO
- **Processing:**          manually annotate 3 times by different persons

  (4 types of labels),

  automatically consolidate,

  automatically and manually assure Annotation Quality

- **Output:**          Json file, containing consolidated annotation and source text url

*Where is ML, still?*

# ANNOTATION WORKFLOW DIAGRAM



- Workflow is cyclic in order to resolve mis-agreement issues
- Lambda functions implement the logic and do powerlifting
- S3 buckets provide structured storage (folders)

*ML is here, finally)*
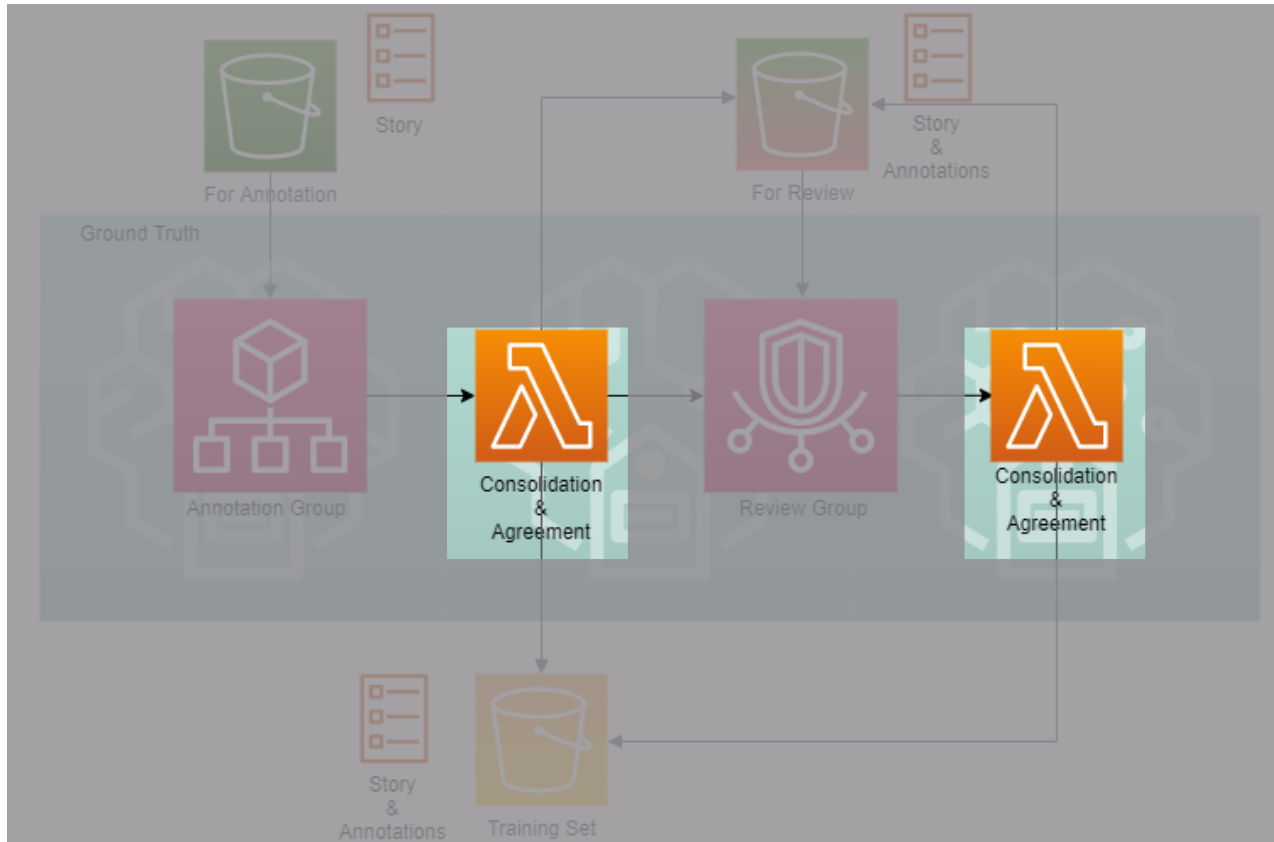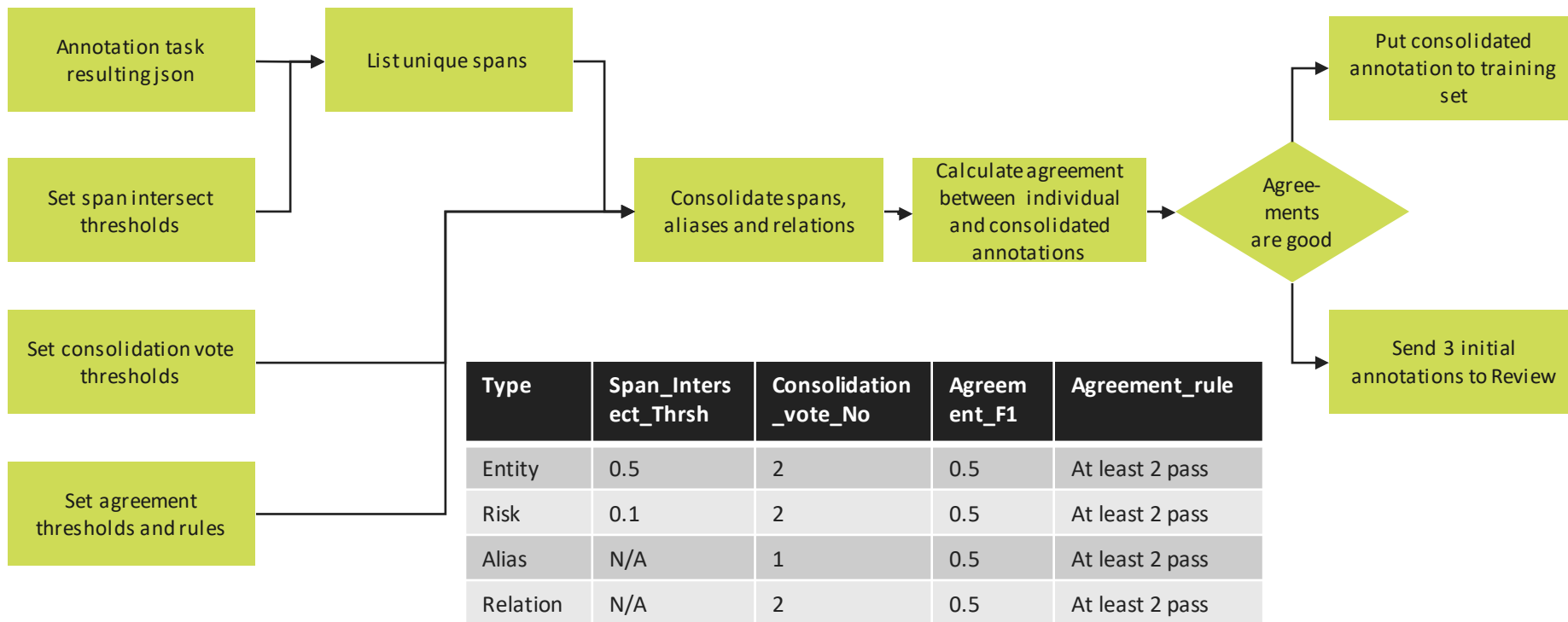
# ANNOTATION TASK



1. Named Entities and their Coreferences
2. Alias Coreferences with NEs
3. Risks
4. Connections

# REVIEW TASK



- Review is made 3 times by SMEs
- If reviews agree, the review replaces initial annotation
- If reviews don't agree, it will be repeated till agreement threshold is passed

# ANNOTATION CONSOLIDATION AND AGREEMENT



- Lambda functions perform the annotation consolidation logic
- They are the gauges of control and adjustment
- End-to-end algorithm performance is the KPI to use when tuning consolidation

# CONSOLIDATION AND AGREEMENT WORKFLOW

Annotation task resulting json

Set span intersect thresholds

Set consolidation vote thresholds

Set agreement thresholds and rules

List unique spans

Consolidate spans, aliases and relations

Calculate agreement between individual and consolidated annotations

Agree-ments are good

Put consolidated annotation to training set

Send 3 initial annotations to Review

| Type | Span_Intersect_Thrsh | Consolidation_vote_No | Agreement_F1 | Agreement_rule |
|------|----------------------|-----------------------|--------------|----------------|
| Entity | 0.5 | 2 | 0.5 | At least 2 pass |
| Risk | 0.1 | 2 | 0.5 | At least 2 pass |
| Alias | N/A | 1 | 0.5 | At least 2 pass |
| Relation | N/A | 2 | 0.5 | At least 2 pass |

# ANNOTATION ADMINISTRATION AND ANALYTICS

## ADMINISTRATION TASKS

- Create users and user groups
- Allocate jobs to specific groups
- Stop and resume annotation tasks
- Find annotation by its attributes
- Show/compare annotations
- Trace consolidated/reviewed annotations to initial ones
- Log annotation sessions (security and time tracking)

## ANALYTICS TASKS

- Annotation velocity: stories, lines, spans per period
- Annotation quality: agreement scores
- Risk topic coverage: snapshot & dynamics
- Annotator's performance: quantity & quality
- Reviewer's performance: quantity & quality
- Training set funnel rates: selected-annotated-consolidated-reviewed-delivered

# QUESTIONS SO FAR?

# WHEN CHOOSE AWS GT: OPEN-SOURCE VS GT

## BRAT

- UI ready for complex text annotation
- Annotation browser + Stop/resume available
- Pairwise comparison of annotated texts
- Security: login, password, IP-address
- Manual user management
- Manual workflow administration
- No annotation analytics

GOOD FIT FOR SMALL-SCALE LEAN TEXT ANNOTATION

## G T

- Automated user management (AWS Cognito)
- Automated workflow (AWS Lambdas, Step functions)
- Automated annotation task fulfillment (count) control
- Security: AWS managed
- UI ready for simple annotation
- No annotation browser, no stop/resume available
- No review/comparison functionality
- No annotation analytics
- Ready for crowed-sourced work force
- Highly scalable

GOOD FIT FOR LARGE SCALE VISUAL ANNOTATION FOR ANN TRAINING

# BRAT UI EXAMPLE: KIND OF OLD-FASHIONED

# GT UI EXAMPLES: CLEAN AND MINIMALISTIC

ADMINISTRATION CONSOLE

Log out

Show instructions

**Jobs**

Start working

< 1 >

| Description | Customer ID | Creation time |
|---|---|---|
| ○ reviewing-Feb-14-a206269-epam-17 | 060725138335 | February 17, 2020 15:12:14 UTC |
| ○ reviewing-Feb-15-a206269-epam-19 | 060725138335 | February 18, 2020 14:34:49 UTC |
| ○ reviewing-Feb-16-a206269-epam-19 | 060725138335 | February 19, 2020 14:04:36 UTC |
| ◉ labeling-Feb-5-a206269- | 060725138335 | February 17, 2020 14:12:18 UTC |

ANNOTATION TASK

Customer ID: 060725138335     Task description: labeling-Feb-5-a20626     Task time: 19:05 of 59 Min 59 Sec     Stop working     Log out

Submit

DataBreach

1 China consumer group raps 2 Apple after 1 ID thefts

UK-APPLE-CHINA 1 China consumer group raps 2 Apple after 2 ID thefts
DataBreach

2018-10-19 10:03:41 GMT+00:00

SHANGHAI (Reuters) - A Chinese consumer body has criticised iPhone maker 2 Apple Inc over a recent data security issue which impacted a
number of consumers in 1 China who said they had 3 suffered financial losses after having their 4 Apple IDs stolen
DataMisuse                                                                                    DataBreach

# REVIEW TASK

# GT CUSTOMIZATION STANDARD OFFERING



- Custom UI template(s)
- Pre-lambda: Task creation
- Post-lambda: Task consolidation

# PROJECT TEAM

- Senior frontend engineer          1 FTE
- Senior .NET/python engineer       1.5 FTE
- Senior cloud devops engineer      1 FTE
- Senior business analyst           1 FTE
- Data scientist                    0.5 FTE

**Project duration – Six month**

# UI CUSTOMIZATION

**REQUIREMENTS:**
- No external libs/dependencies
- Unicode-compatible
- Works in latest Chrome and Firefox

**TECHNOLOGY:**
- HTML
- CSS Grid Layout
- SVG
- JS (ES6)
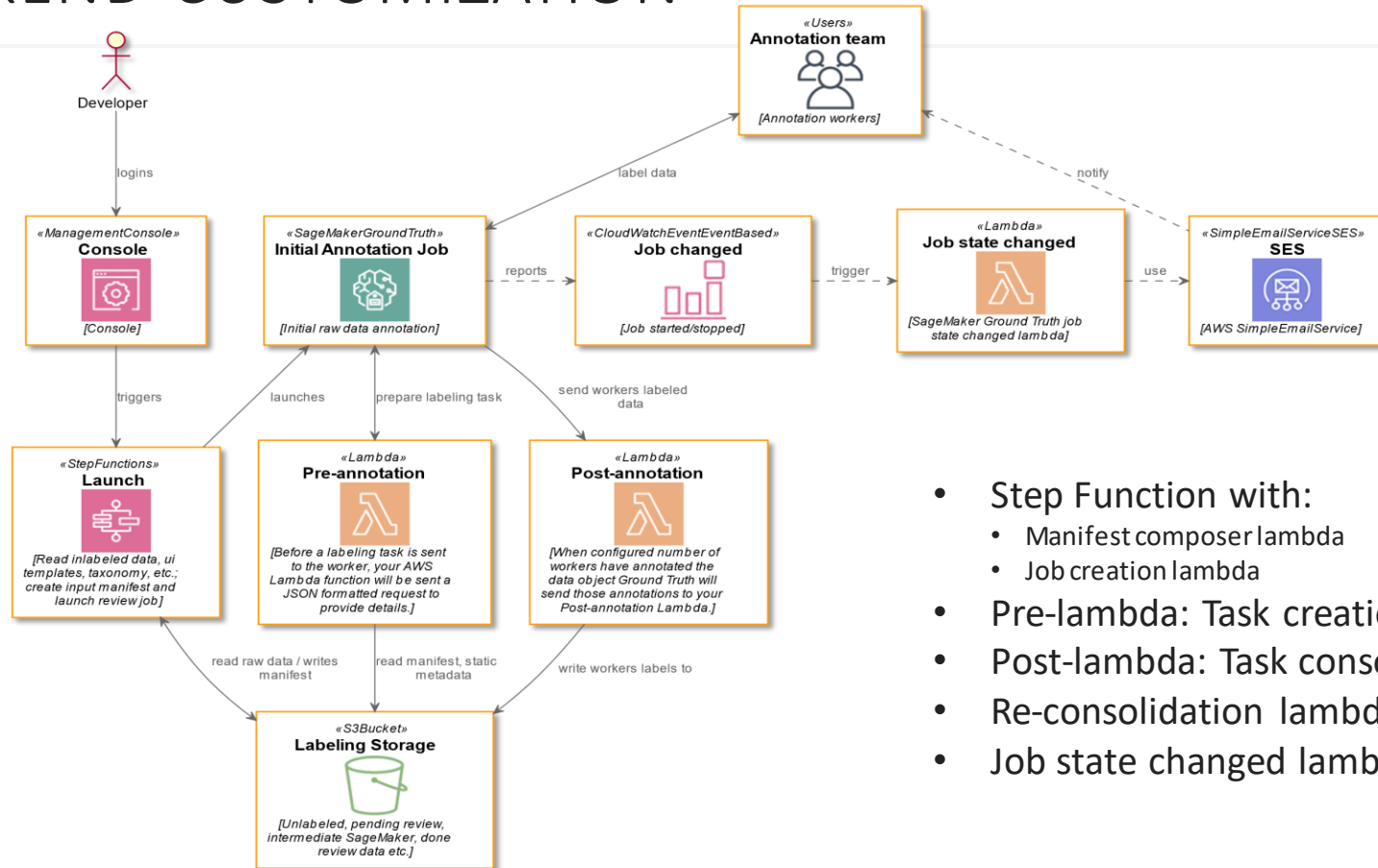- HTML DOM selection API

| CHALLENGE | SOLUTION |
|---|---|
| JS string functions, DOM selection API and Python backend count characters differently | Custom character counting function based on ES6 features; Unicode friendly |
| Visualize relations between spans as arches in html | Dynamically generated in-line SVG |
| DOM selection API is quite stubborn one | Patience) |

```
      path"></path>
    <path d="
    M 69.640625,741.984375
    Q 42.953125,661.984375
    16.265625,1114.984375" class="entity1-
    path"></path>
  </svg>
▼<div class="text-container">
  ▼<h1>
      "Trump signs order aimed at stopping "
    ▼<span class="risk risk1" data-num="1"
    id="lc7c4689risk1" data-type=
    "IllicitGoods" data-type-id="21">
        ::before
        "sale of counterfeit products" == $0
        ::after
    </span>
    " through "
    ▼<span class="entity entity1 super-l1"
    data-num="1" id="lc7c4689entity1">
        ::before
      ▼<span class="entity entity2" data-
      num="2" id="lc7c4689entity4">
          ::before
          "Alibaba"
      </span>
```

# BACKEND CUSTOMIZATION



- Step Function with:
  - Manifest composer lambda
  - Job creation lambda
- Pre-lambda: Task creation
- Post-lambda: Task consolidation
- Re-consolidation lambda
- Job state changed lambda

# S3 OPERATIONAL STORAGE



STORY FILES ARE MOVING FROM FOLDER TO FOLDER AS THEY PROGRES ALONG PIPELINE TO "GOLDEN DATA SET"

Folders:
- Documents_xml
- Annotations_inprogress
- Review_selected
- Review_inprogress
- Annotations_done

# ANALYTICAL DWH: DYNAMO vs ATHENA vs AURORA



| | Dynamo | Athena | Aurora |
|---|---|---|---|
| Analytics fitness | ☹ | ☺ | ☺ |
| Implementation complexity | ☺ | ☹ | ☺ |
| Cost of ownership | ☺ | ☹ | ☺ |

# ANOTATION KNOWLEDGE MANAGEMENT

- Annotation rules
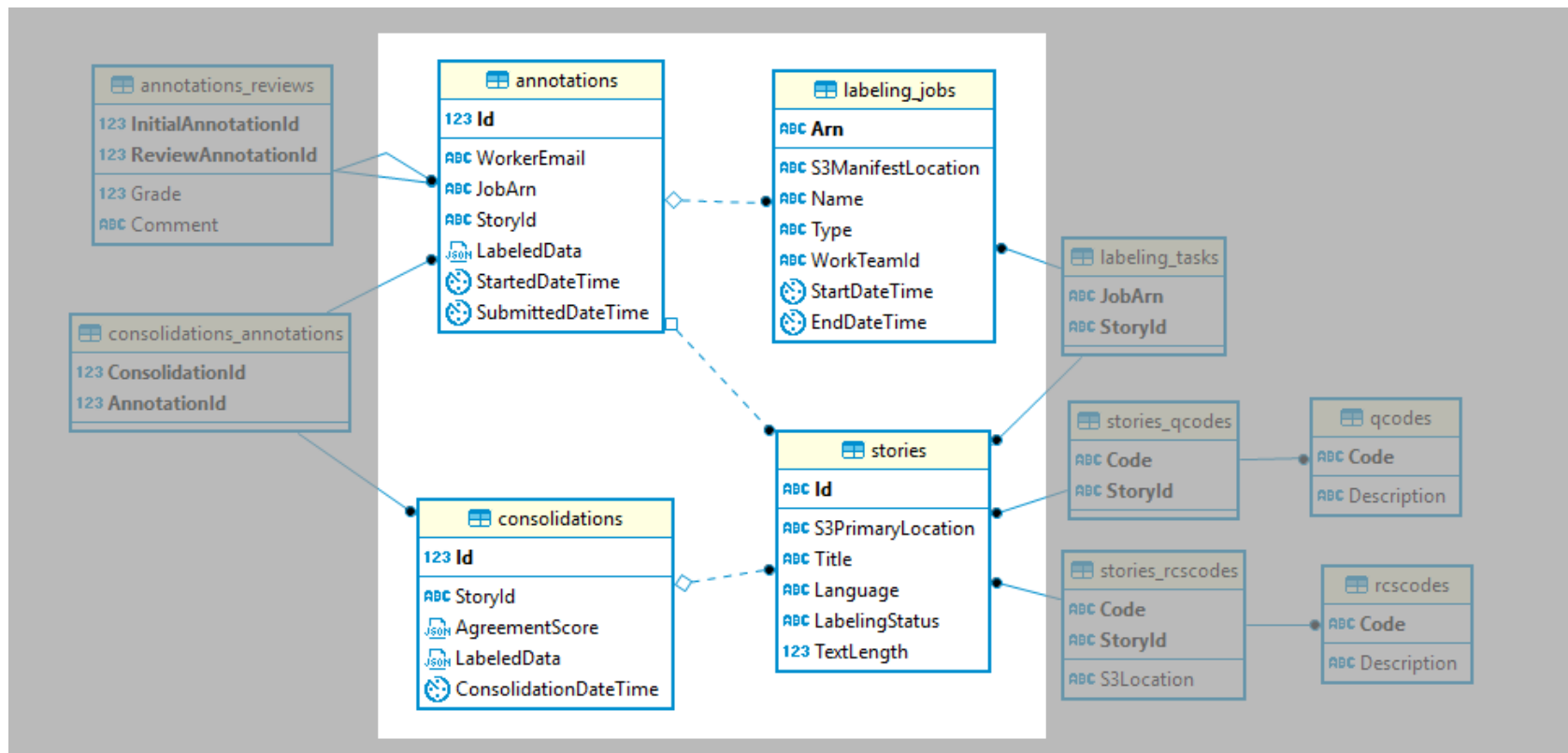- Annotation FAQ
- Crime taxonomy 2.3
- Corporate crime :  Guidance on how tagging should be applied for this part of the taxonomy
- Non-corporate crime: Guidance on how tagging should be applied to this part of the taxonomy
- Cases for annotation quality improvement

CORRUPTION VS BRIBERY
DRUG TRAFFICKING VS DRUG ABUSE
SEXUAL EXPLOITATION VS SEX OFFENCE
THEFT VS EMBEZZLEMENT VS ROBBERY VS BURGLARY
GENERAL FRAUD VS SECURITIES FRAUD

"DUMMY" NAMED ENTITIES
SANCTIONED ENTITIES
TITLES & OCCUPATIONS
SUPER- AND SUB-ENTITIES
GENERIC CRIMES

# QUESTIONS?

# ANOTATION OUTCOME STATS

- Annotation team size:             5 persons
- Annotation period:                 6 month
- Stories annotated:                 1200
- Risk coverage min-avg-max :    1-50-420
- Stories reviewed:                  ~ 50
- Median NE agreement:           0.92
- Median alias agreement:         0.85
- Median risk agreement:          0.98
- Median relation agreement:      0.86

# LESSONS LEARNED

- Comprehensive annotation tool is complex and expensive

- Ground truth is difficult to adopt to complex text annotation

- **Annotation tool is only a tool**

- Good taxonomy is very important

- Annotation rules body of knowledge is a must

- Expect lot of investments into training, review and re-training

# FINAL JUDGEMENT

- Tool is used by several teams within the customer organization

- Nice corporate toy to brag about

- **Proof of investment is difficult, maybe, better spend this money for more qualified and motivated annotators?!**

Whoops! Demolition company accidentally tears down wrong house.

Burglary OR Theft?? 😊