

# SQUARE: A Benchmark for Research on Computing Crowd Consensus

**Presentation by**  
**Srikanth Muralidharan**

# Case Study – Brexit

- Referendum to leave/stay
- 17.4 M wanted to leave
- 16.1 M wanted to stay back
- Result: By Majority vote, 'leave' by margin of 1.3 M

[Source: Independent news article](#)

# Case Study – Brexit



INDEPENDENT

News

Voices

Culture

Lifestyle

Tech

Sport

Daily Edition

News › UK › UK Politics

## Brexit research suggests 1.2 million Leave voters regret their choice in reversal that could change result

The research suggests that if a second referendum were held, the vote would be much closer

Lizzie Dearden | @lizziedearden | Friday 1 July 2016 |  735 comments

[Source: Independent news article](#)

**REMEMBER MARGIN WAS 1.3 M!**

# Case Study – Brexit

**TOP QUESTIONS ON THE EUROPEAN UNION** Google Trends  
in the UK since Brexit result officially announced

- 1 What does it mean to leave the EU?
- 2 What is the EU?
- 3 Which countries are in the EU?
- 4 What will happen now we've left the EU?
- 5 How many countries are in the EU?

[google.com/trends](https://google.com/trends)



**GoogleTrends** ✓  
@GoogleTrends

[Follow](#)

"What is the EU?" is the second top UK question on the EU since the [#EURefResults](#) were officially announced

4:25 AM - 24 Jun 2016

↩ ↻ 27,217 ❤ 19,072

[Source: Fortune](#)

**MAJORITY VOTE CONSENSUS DOESN'T SEEM OPTIMAL HERE!**

# Take Away's

- Getting Accurate Data label is a challenging process
- Possible solution: Get labels from multiple workers for a given data, combine them

# Take Away's

- MV might not be optimal method for combining.
- Better consensus methods might consider worker (voter) expertise / reliability and so on.

# Problem

- Multiple consensus methods exist in literature
- Compare and contrast between different consensus methods
- Design rules about when to use a consensus method

# Contributions

- SQUARE Benchmark
  - “(Statistical Quality Assurance Robustness Evaluation”
- Collection of diverse datasets, consensus methods
- Simulation of varying noise, supervision
- Heuristics for picking right consensus method given the scenario



# Outline

- Benchmark Datasets for comparisons
- Consensus methods
- Experiment Setup and Results
- Concluding Remarks

- **Benchmark Datasets for comparisons**
- Consensus methods
- Experiment Setup and Results
- Concluding Remarks

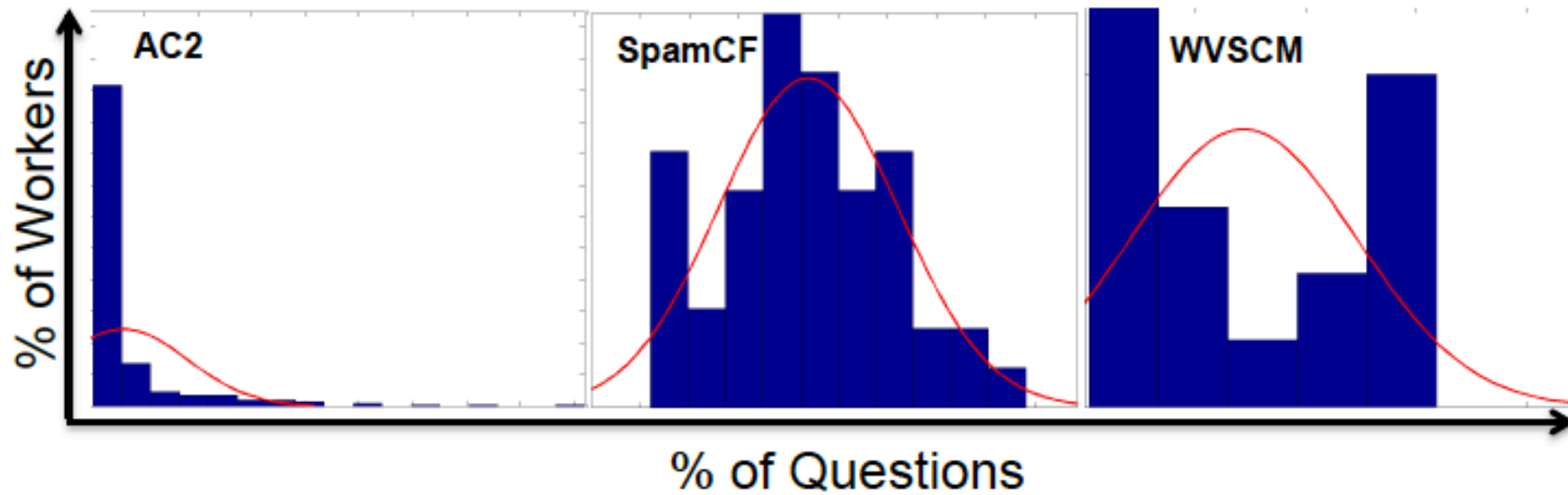
# Datasets

<b>Dataset</b>	<b>Categories</b>	<b>Examples</b>	<b>Workers</b>	<b>Labels</b>	<b>MV Acc.</b>
AC2	4	333	269	3317	88.1
BM	2	1000	83	5000	69.6
HC	3	3275	722	18479	64.9
HCB	2	3275	722	18479	64.8
RTE	2	800	164	8000	91.9
SpamCF	2	100	150	2297	66.0
TEMP	2	462	76	4620	93.9
WB	2	108	39	4212	75.9
WSD	3	177	34	1770	99.6
WVSCM	2	159	17	1221	72.3

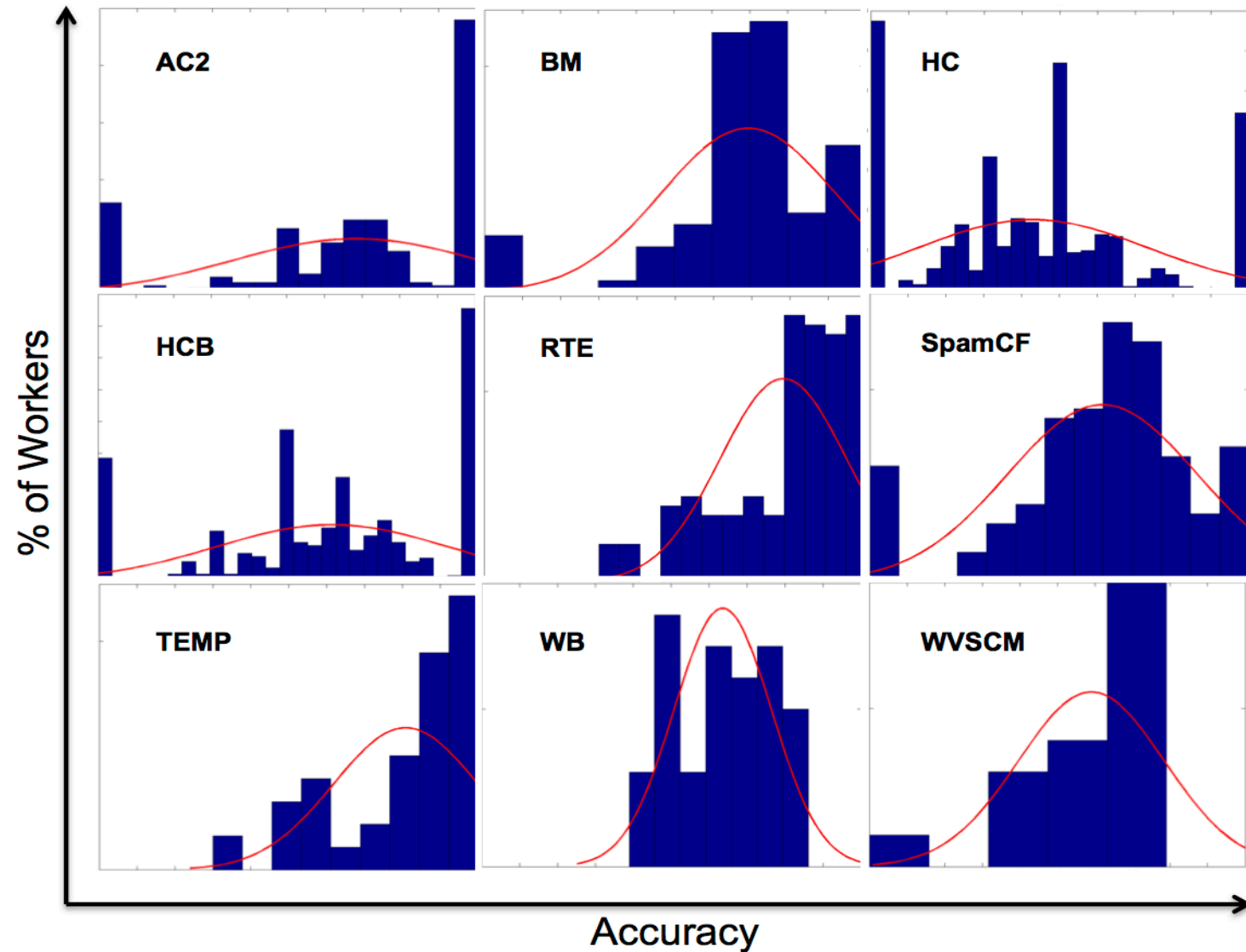
# Datasets

<b>Dataset</b>	<b>Categories</b>	<b>Examples</b>	<b>Workers</b>	<b>Labels</b>	<b>MV Acc.</b>
AC2	4	333	269	3317	88.1
BM	2	1000	83	5000	69.6
HC	3	3275	722	18479	64.9
HCB	2	3275	722	18479	64.8
RTE	2	800	164	8000	91.9
SpamCF	2	100	150	2297	66.0
TEMP	2	462	76	4620	93.9
WB	2	108	39	4212	75.9
WSD	3	177	34	1770	99.6
WVSCM	2	159	17	1221	72.3

# Labels per worker count

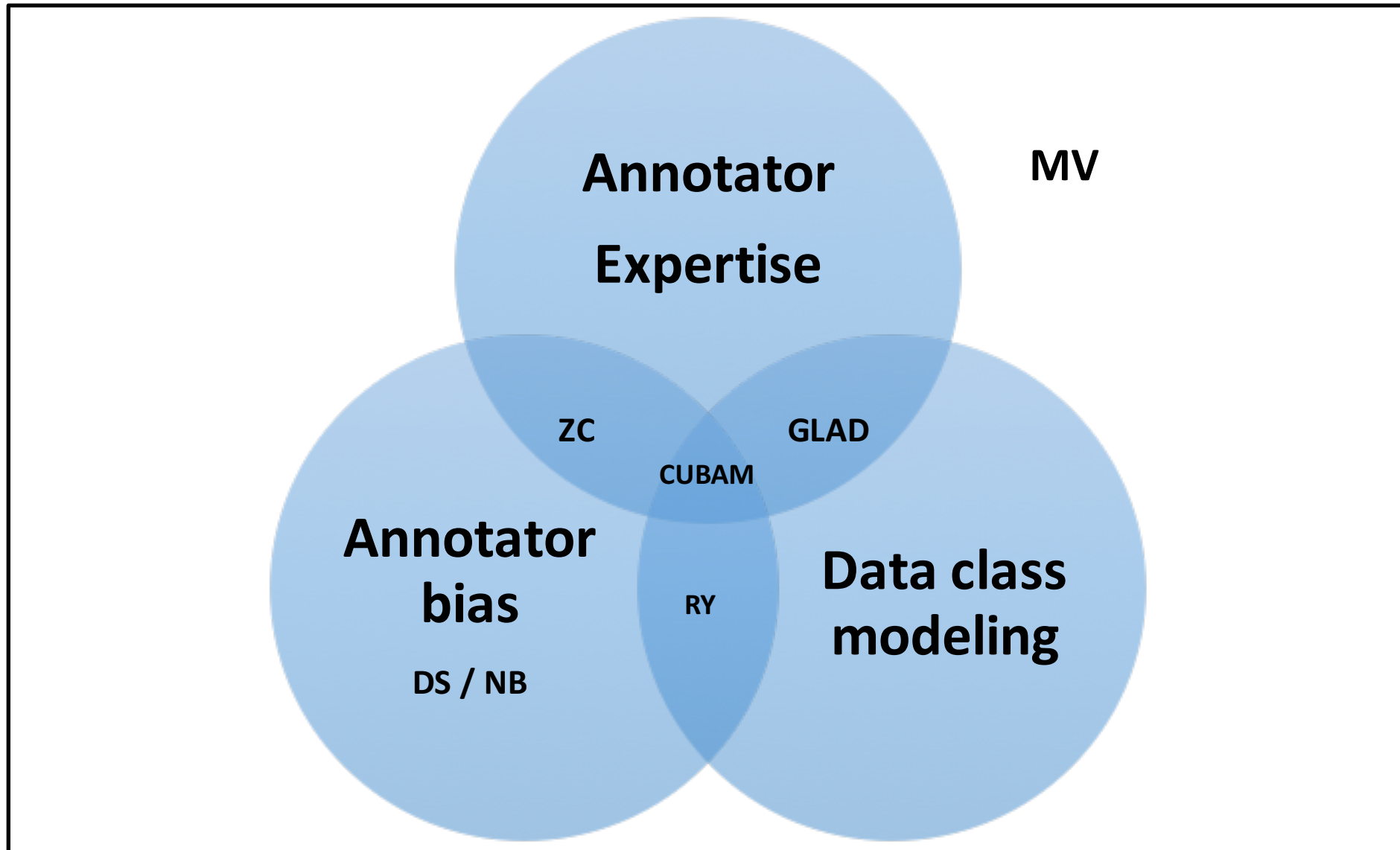


# Worker accuracy distribution



- Benchmark Datasets for comparisons
- **Consensus methods**
- Experiment Setup and Results
- Concluding Remarks

# Modelling considerations





# Consensus methods – An overview

- Many methods exist for consensus
- Vast variations present between each methods

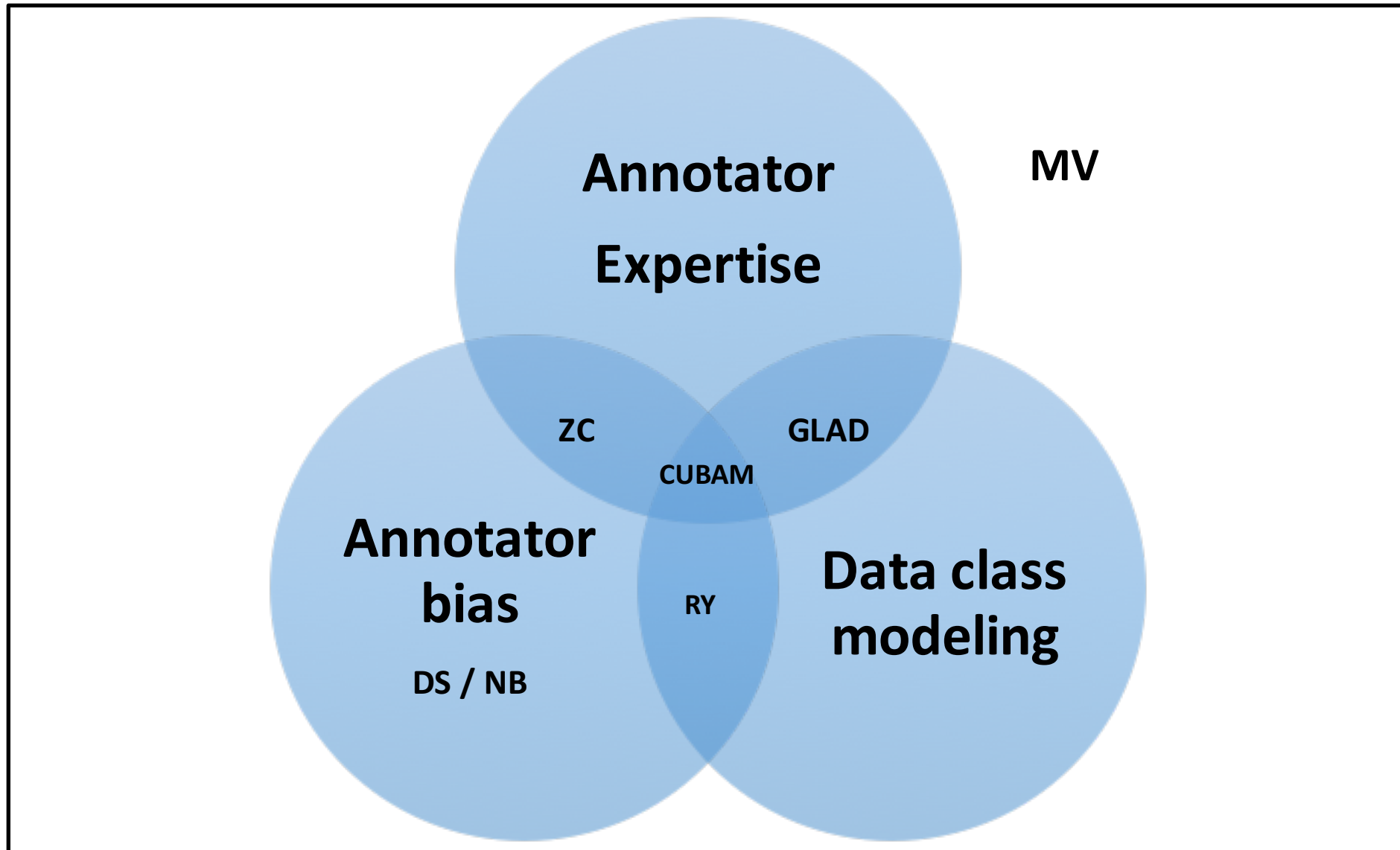
# Key sources of Diversity

- Modelling of Worker Behaviour
- Modelling of class label distribution

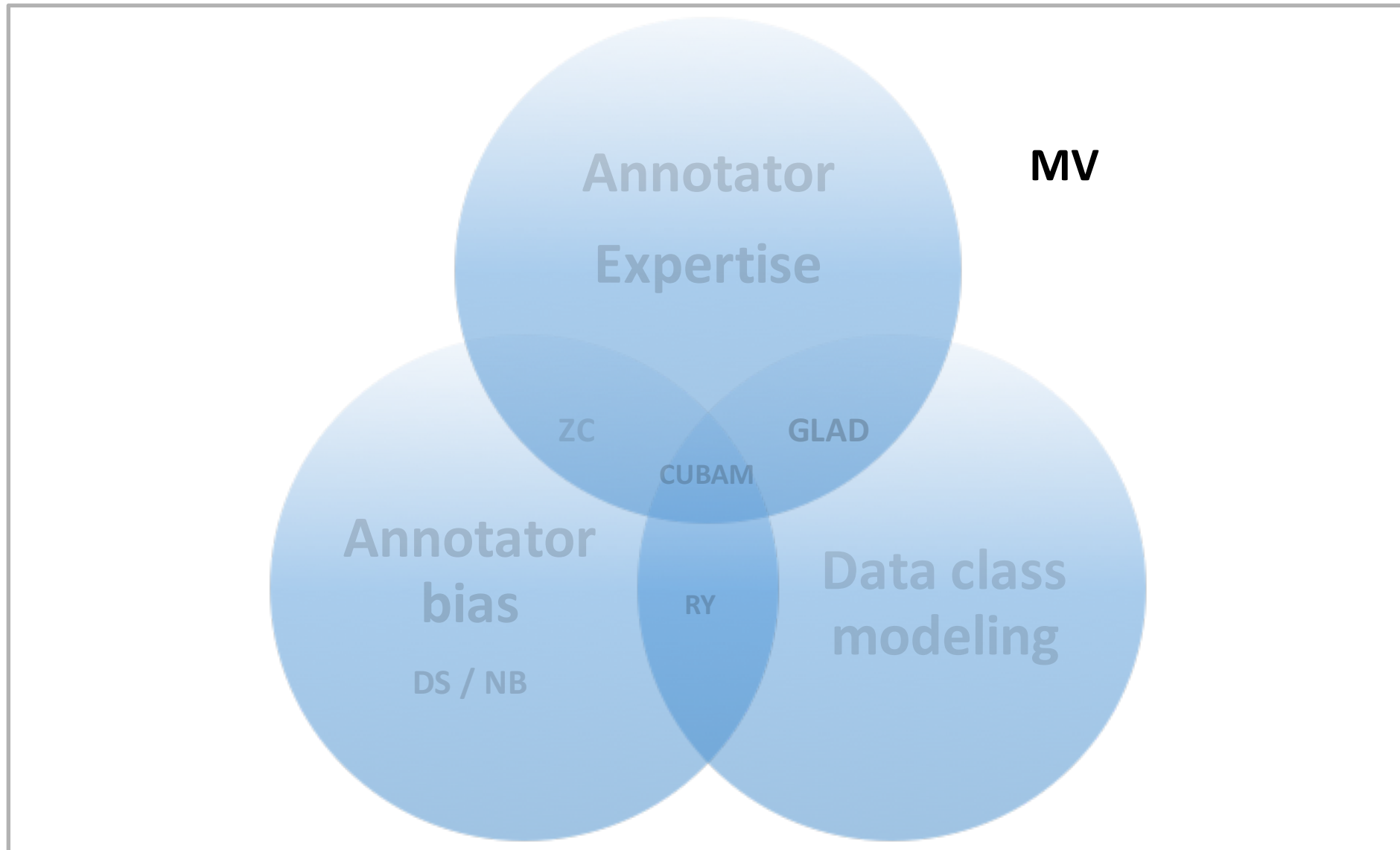
# Consensus methods used in our experiments

- Majority voting (MV)
- ZenCrowd (ZC) [Demartini et al. 2012]
- David and Skene (DS) & Naïve Bayes (NB)
- GLAD [Whitehill et al. 2009]
- RY [Raykar et al. 2010]
- CUBAM [Welinder et al. 2010]

# Modelling considerations



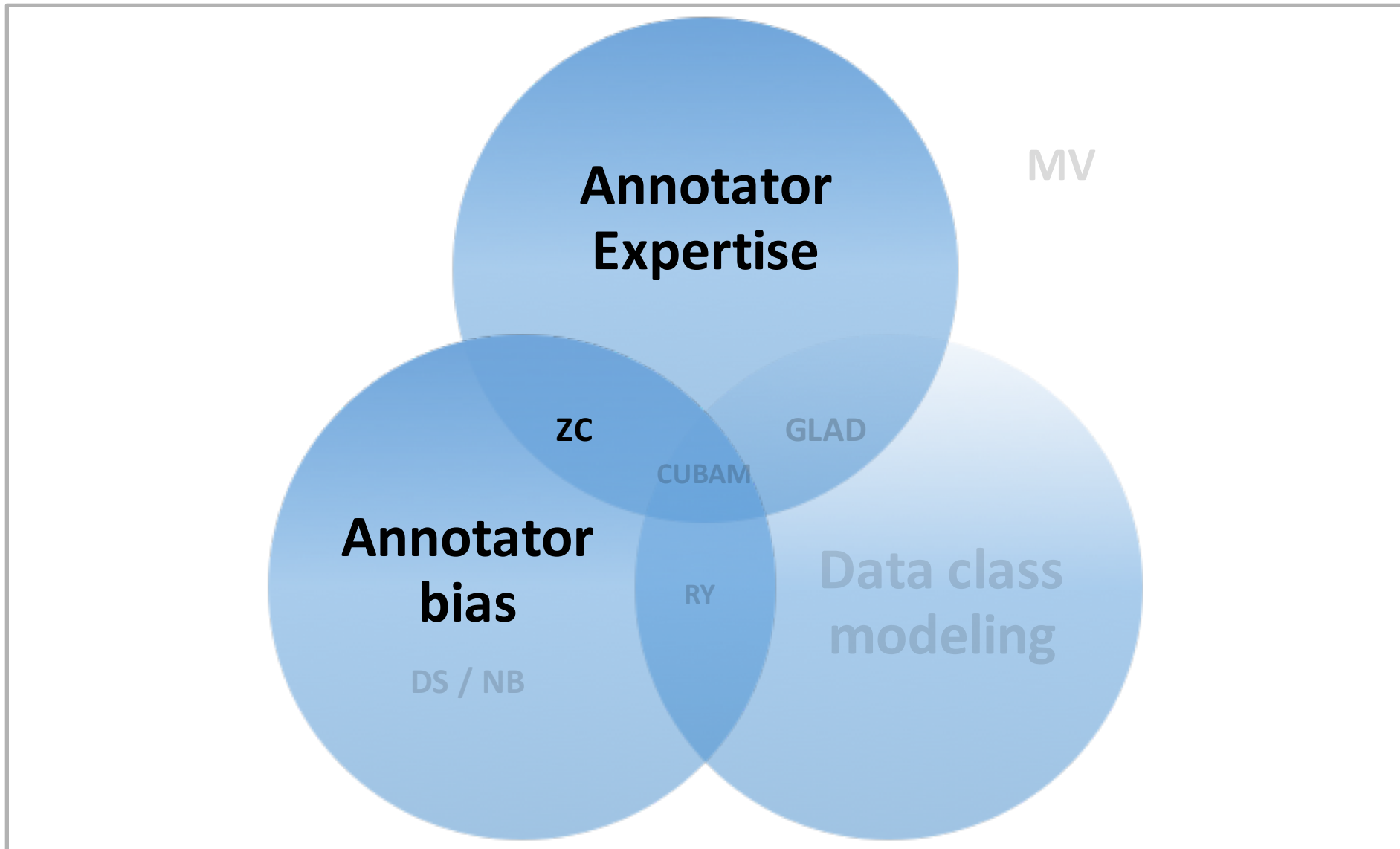
# Modelling considerations



# Majority voting

- Most widely used consensus method
- Assumes workers to be high quality and Independent, imply
  1. No modeling of worker behavior
  2. No modeling of annotation process
  3. Independent of task
- No estimations -> Fast label inferencing
- Applied under all supervision-classification modes

# Modelling considerations



# ZenCrowd

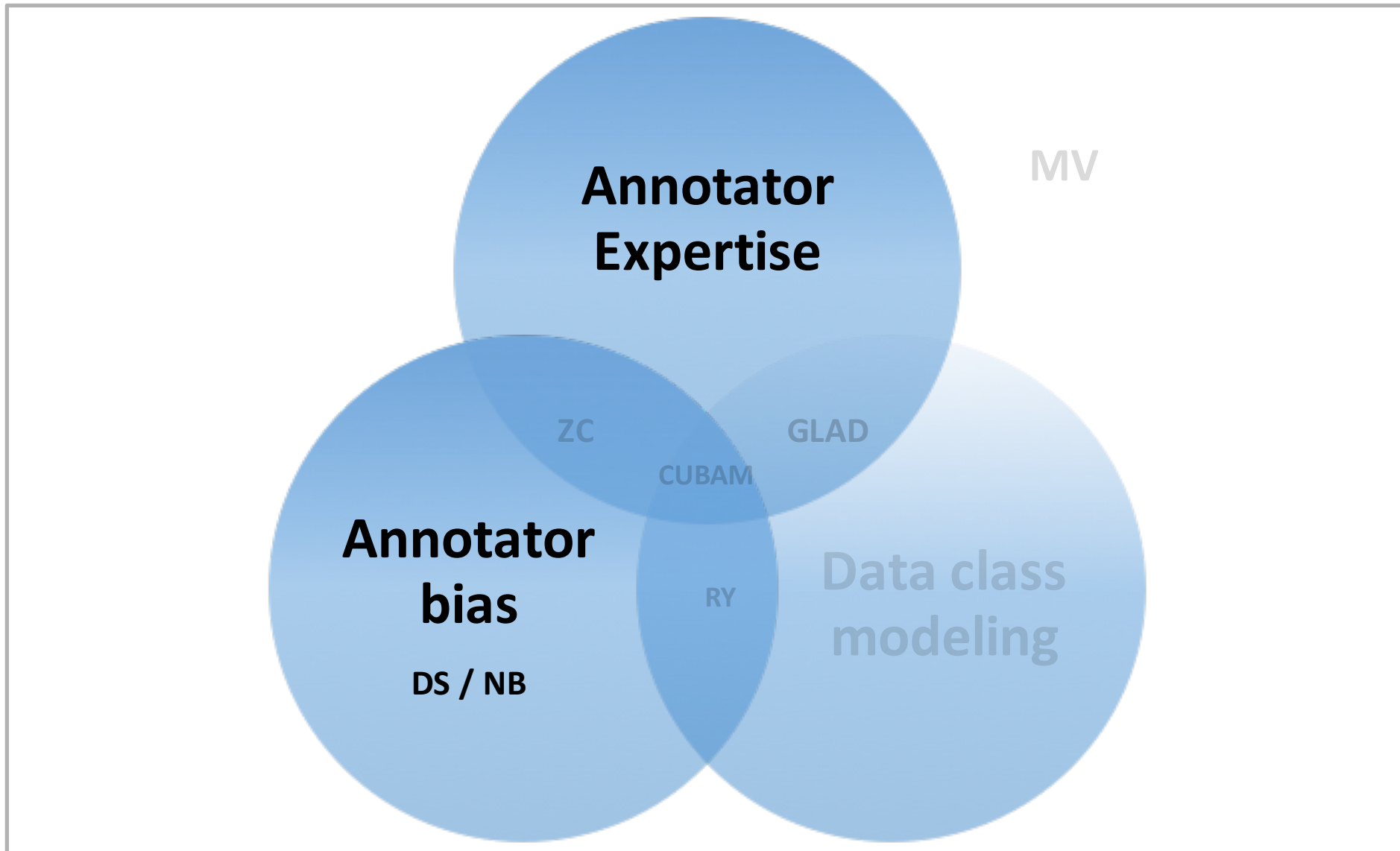
- Parameterizes worker reliability
  - Single Parameter per worker
  - Measure of worker expertise
  - Easily Generalizable to multi-class / Multi choice tasks
  - Easy mechanism to identify and handle Adversarial workers (bias)
- Other assumptions apply as in MV



# More about ZenCrowd

- Unsupervised in original version
  - Use Expectation-Maximization (EM) to estimate worker reliability and labels
- Could be operated under other supervision modes
  - Lightly Supervised by only providing class priors
  - Semi Supervised by providing gold labels for a subset [Wang et al. 2011]
  - Fully Supervised using MLE [Snow et al. 2008]

# Modelling considerations



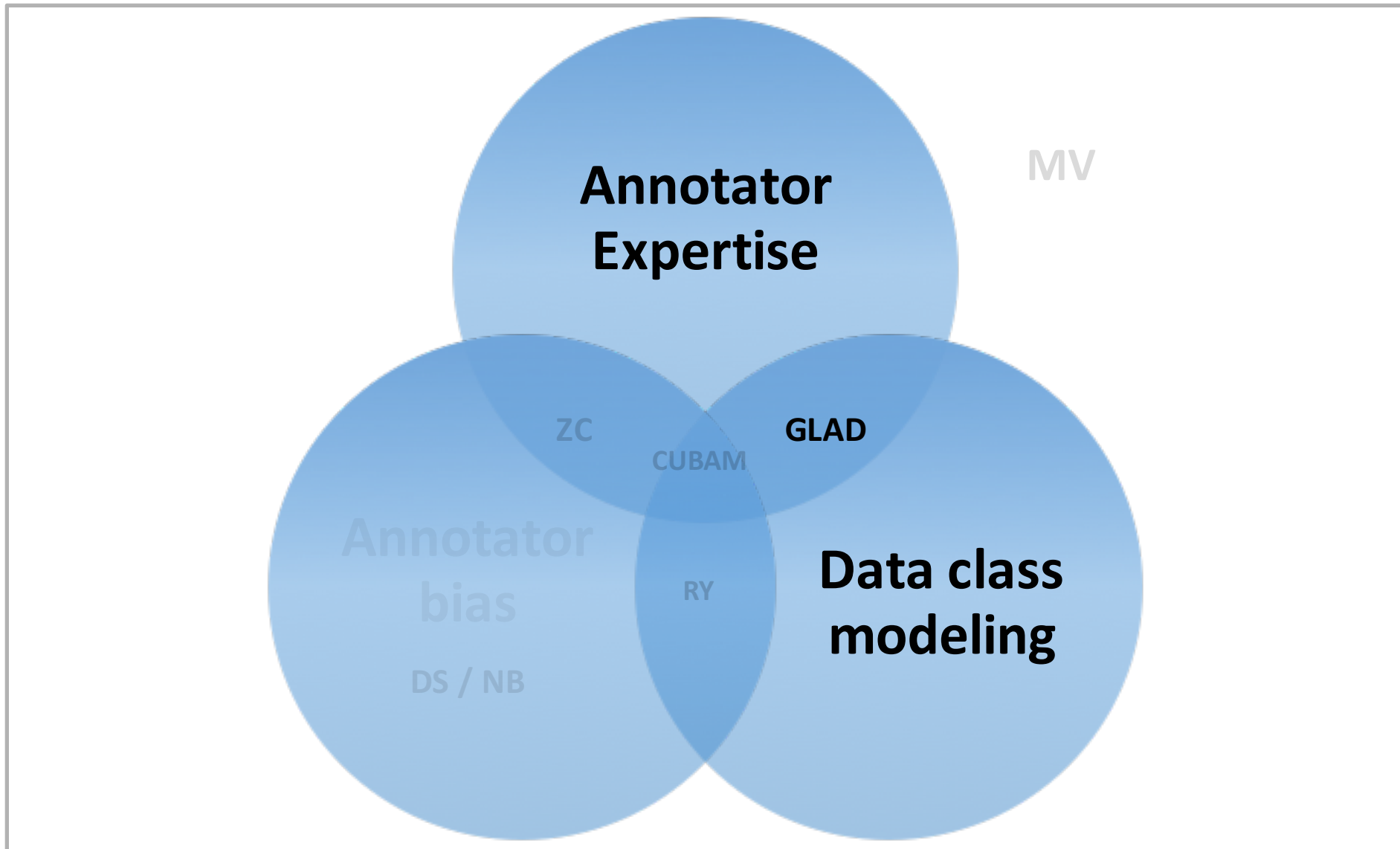
# Dawid Skene & Naïve Bayes

- Parameterizes worker bias w.r.t all classes
- Models Confusion Matrix for each worker and class prior
  - Implies one parameter per class per worker
  - Easy detection of perfections/imperfections of each worker w.r.t a class
  - Sparsity could be problematic.
  - Cannot be generalized to multiple choice Selection task
- Other assumptions apply as in MV

# More about DS / NB

- Unsupervised in original version
  - Use EM to estimate worker CM and class prior
- Could be operated under other supervision modes
  - Under Lightly Supervised and semi-supervised conditions
    - Variant Estimation technique to distinguish b/w noise and worker bias
  - Under Full supervision by MLE using Laplacian smoothing

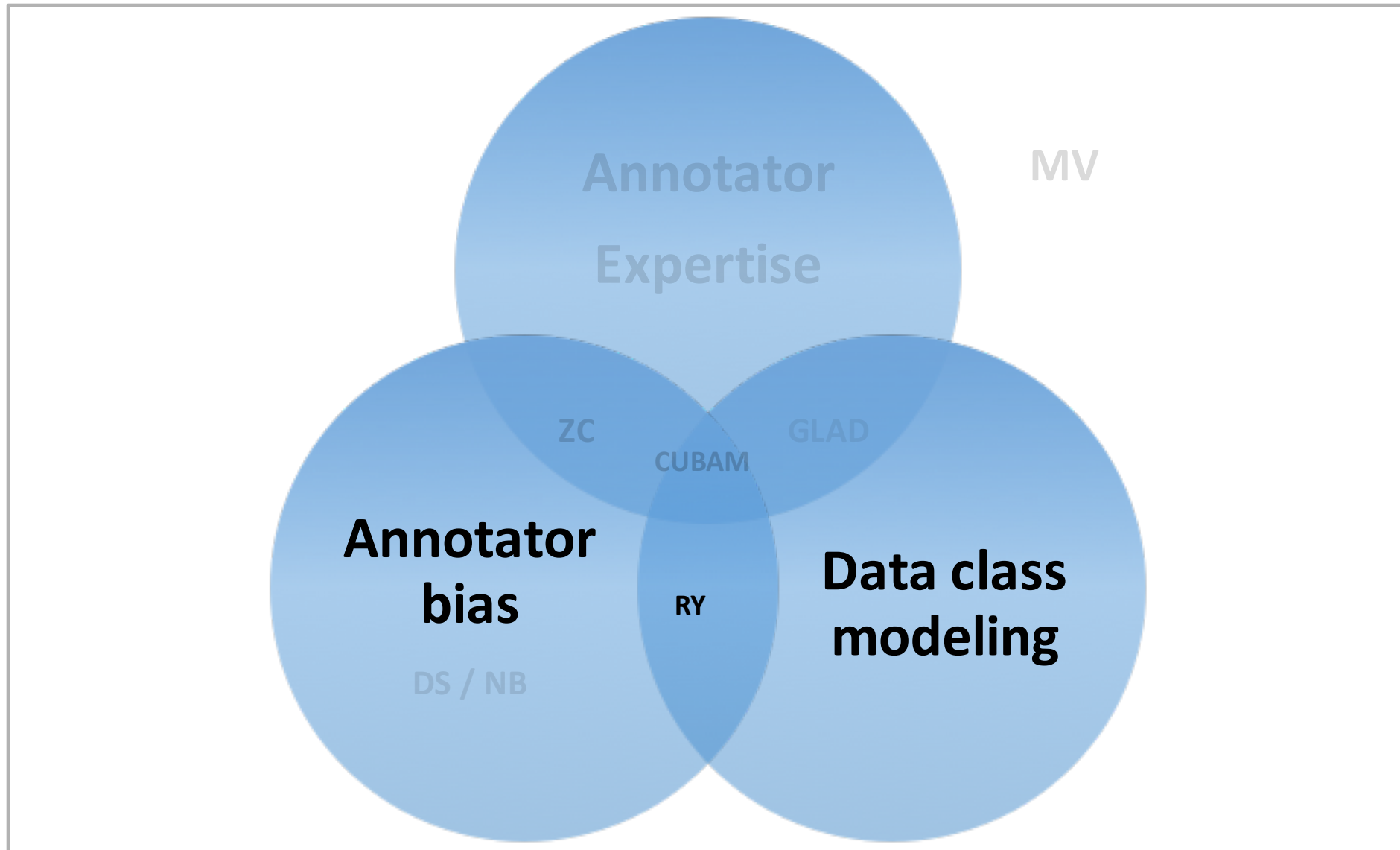
# Modelling considerations



# GLAD

- Joint Parameterization of worker expertise and class label (though example difficulty modeling)
- Originally proposed as unsupervised problem
  - Addressed using EM with gradient ascent M step
  - Extendable to Multi-class problems
  - Could be projected under other supervised settings (Fit labels in EM)

# Modelling considerations

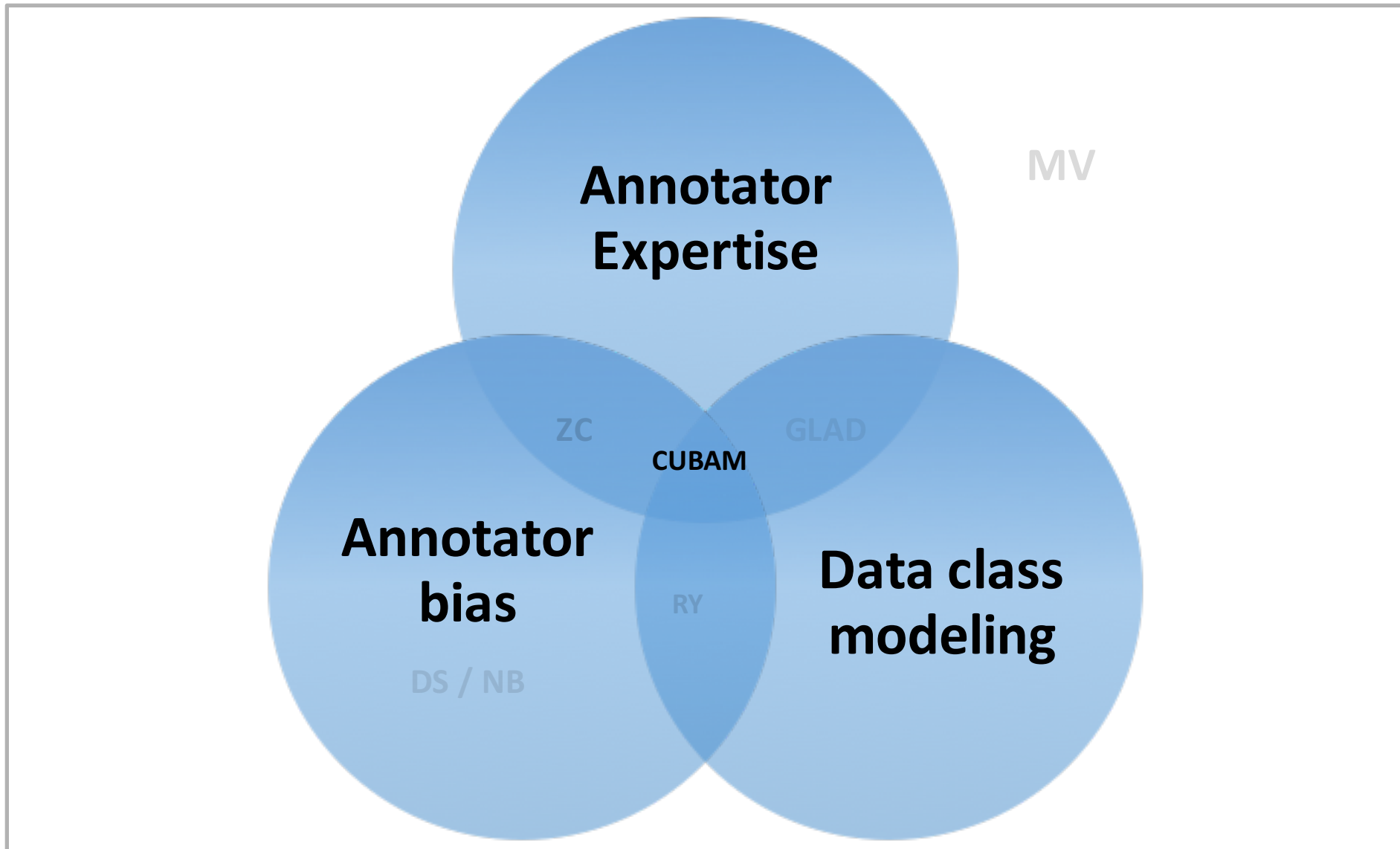


# RY (Raykar)

- Joint parametrization of worker bias for each class
  - Specificity and Sensitivity in binary case modeled using beta prior
  - Dirichlet prior in multi-class case
- Optional: Feature representation of example
  - Used to estimate labels if available, DS like mechanism otherwise
- Originally proposed as unsupervised problem
  - EM to estimate worker bias parameters
  - Extendable to Multi-class problems
  - Class (example) parameters => No multi-choice



# Modelling considerations



# CUBAM

- Incorporate all the three modalities
  - Normalized relevance weight to each worker
  - Worker labels could be determined, given input and worker specific parameters
  - Labels obtained by MAP estimation of worker specific parameters and input
  - Worker expertise/bias could be determined.
- Multi-class classification possible
- No direct supervision is apparent

# Summary

	Unsup	Light sup.	Semi sup.	Full sup.	Multi-class	Mult-choice Q
MV	✓	✓	✓	✓	✓	✓
ZC	✓	✓	✓	✓	✓	✓
DS	✓	✓	✓	✓	✓	✗
NB	✗	✗	✗	✓	✓	✗
GLAD	✓	✓	✓	✓	✓	✓
RY	✓	✓	✓	✓	✓	✗
CUBAM	✓	✗	✗	✗	✓	✗

- Benchmark Datasets for comparisons
- Consensus methods
- **Experiment Setup and Results**
- Concluding Remarks

# Goal

- Project consensus under exhaustive set of settings
- Analyze importance of each assumption
- Analyze impact of different conditions on the performance (e.g. supervision)

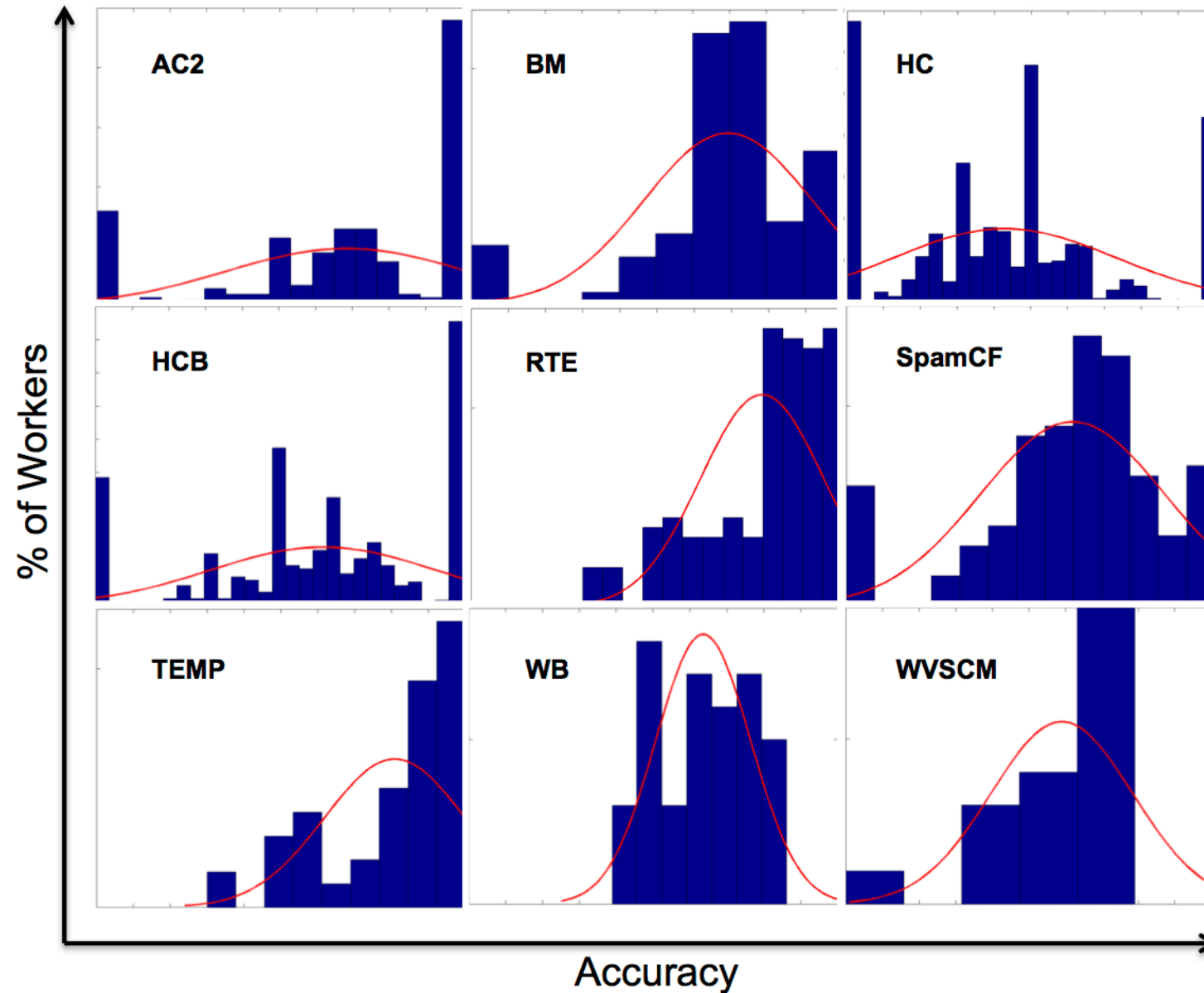
# Possible settings

- Dataset
- Degree of Supervision
- Label Noise
- Evaluation Metrics

# Degree of Supervision

- Ranging from unsupervised to 90% supervised
- Cross-validation employed in each case, with varying folds
- Unsupervised – No information about hyperparameters/priors
- Light supervision – Data Prior distribution available
- Full Supervision – Gold-labeled training examples available

# Sampling visualization





# Label Noise

- Noise added to worker label predictions
- Noise added under maximizing data realism
  - Preservation of worker accuracy distribution
  - Preservation of worker-example correspondences
- Worker accuracy parameterized by normal distribution
- Different proportions of noise added
  - For each worker, accuracy sampling is performed, new labels generated

# Evaluation Metrics

- Ideal metric: Simplify understanding, easy adaption
- F1 and accuracy metrics employed here
- Other alternatives possible – e.g. Significance testing [Smucker et al. 2007]

# Details about consensus implementations - ZC

- Beta priors used for worker reliability
- Dirichlet priors used for class label distribution
- In unsupervised setting, uniform distribution was used for category distribution modeling
- Parameters computed using training data in supervised setting

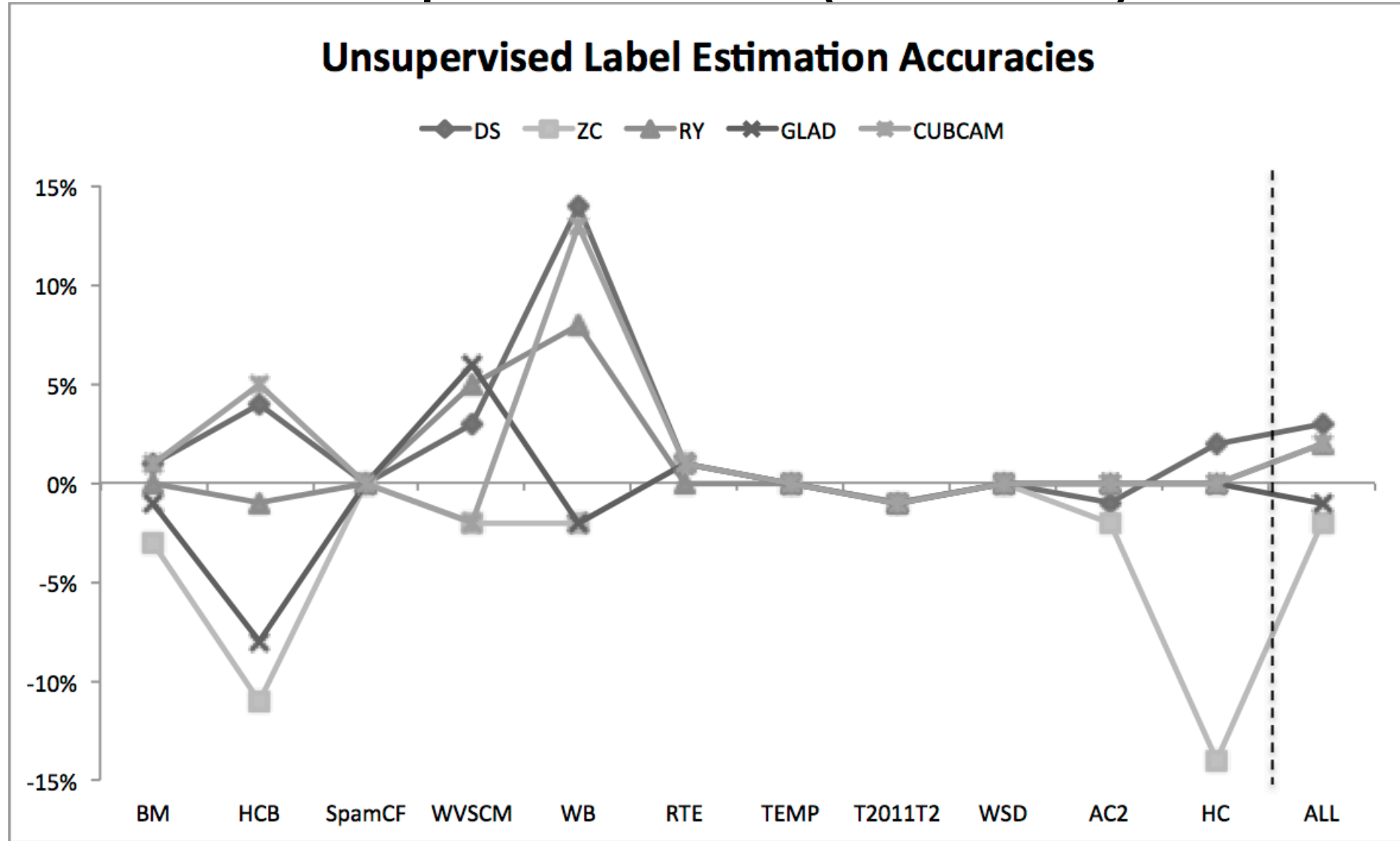
# Details about consensus implementations - RY

- Beta priors used for worker specificity – sensitivity
- Only used for binary classification as in original paper [Rayker et al. 2010]
- In unsupervised setting, parameters were assumed as in ZC
- Parameters computed using training data in supervised setting

# Details about consensus implementations – CUBAM, DS, GLAD

- CUBAM – same prior assumptions as original implementation
- DS – priors not assumed – computed using training data when available
- GLAD – uniform class distribution assumed – computed using training data when available
- CUBAM / GLAD – support only binary classes

# Results on unmodified datasets – Unsupervised (vs MV)



# Results on unmodified datasets – Absolute Numbers

[illegible]

# General Trends – unmodified datasets

- DS – top average performance under  $< 50\%$  supervision
- RY – top average performance under  $> 70\%$  supervision
- CUBAM – performs well under noisy conditions of HCB
- MV - comparable to others under noisy and scarce data

[illegible]



# General Trends – unmodified datasets

- Performance limited by worker/class prior modeling
- Smart modelling of workers => Good performance
  - Sufficiency in model complexity further gets aided by supervision

# Results on datasets – Different Noise injections

[illegible]

# General Trends – Noisy datasets

- MV – Outperformed at modest noise levels (50%)
- DS – top average performance under low noise
- CUBAM – doesn't perform better than others
- RY – performs well with supervisions – highlights significance of priors

[illegible]

# General Trends – Overall

- MV – Outperformed when noise levels are high
- DS – top average performance under low noise
- RY – top average performance under high supervision – highlights significance of domain knowledge
- CUBAM – doesn't perform better than others – shows performance sensitivity to datasets
- ZC and GLAD perform similarly

- Benchmark Datasets for comparisons
- Consensus methods
- Experiment Setup and Results
- **Concluding Remarks**


# Conclusion




- Design of a comprehensive benchmark for evaluating consensus methods
- Diverse dataset – consensus methods to rely on
- Diverse experiments to compare and contrast different consensus methods

# Future work

- Improved Benchmark tests for better analysis
- Tuning of current methods for fair comparisons
- Analysis under sampling from worker empirical distribution rather than normal distribution

# SQUARE API on web

 [utir / square](#)

 Watch 2  Star 4  Fork 3

[Code](#) [Issues 2](#) [Pull requests 0](#) [Projects 0](#) [Pulse](#) [Graphs](#)

SQUARE (Statistical QUALity Assurance Robustness Evaluation) <http://ir.ischool.utexas.edu/square/>


14 commits








1 branch


0 releases

2 contributors

Branch: master [New pull request](#) [Find file](#) [Clone or download](#)

 **aashish-sheshadri** Update readme.md Latest commit effb036 on Oct 27, 2015

 <a href="#">processingScripts</a>	Optimize imports	a year ago
 <a href="#">sampleCLA</a>	Adding support to aggregation without ground truth and a Matlab scrip...	3 years ago
 <a href="#">src</a>	Use generics to resolve compiler warnings	a year ago
 <a href="#">.gitignore</a>	Update dependencies and enforce better builds	a year ago
 <a href="#">LICENSE.md</a>	Create LICENSE.md	3 years ago
 <a href="#">pom.xml</a>	Add -Xlint:unchecked to pom.xml	a year ago
 <a href="#">readme.md</a>	Update readme.md	11 months ago

 **readme.md**

**October 26, 2015: Version 2.0 of SQUARE now released! (new GIT repo to avoid impacting those using Version 1.0). See 2.0 page for details of what's new! -- [SQUARE-2.0](#).**

Build and Run Instructions: -- Project conforms to a maven build. -- Run mvn install from the root directory. -- See \_CLA files for example usage. See: <http://ir.ischool.utexas.edu/square/> for results analysis and API.

Processing scripts: -- This folder includes Matlab scripts to interface with data generated by the SQUARE algorithms to plug and play with other algorithms. — Generating data 1. Create a folder of your choice 2. The folder needs to contain the files: a)categories.txt - each line names the category b)responses.txt - space separated values with each line of the format workerId question response c)groundTruth.txt - space separated values with each line of the format question response Note that its not necessary to have the groundTruth.txt file 2. Depending on whether you have ground truth or



Thank you