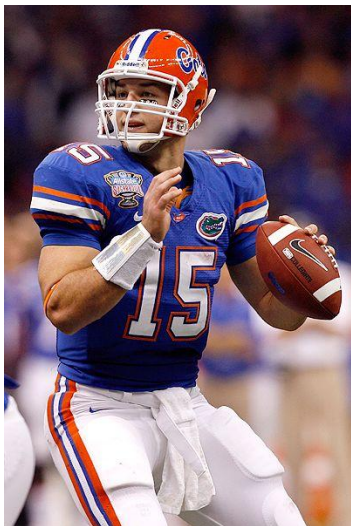


# University of Florida DSR Lab System for KBP Slot Filler Validation 2015

---

Miguel Rodriguez, Sean Goldberg, Daisy Wang

# Slot Filler Validation



Tim Tebow

gpe:schools\_atended



Bristol Central High School
New England Patriots
University of Florida
University of Connecticut
ABC News

# Slot Filler Validation



Tim Tebow

gpe:schools\_atended



	Truth
Bristol Central High School	T
New England Patriots	F
University of Florida	T
University of Connecticut	F
ABC News	F

# Slot Filler Validation



org:subsidiaries



	Truth
Survey Research Center	T
Florida Museum of Natural History	T
Smithsonian Tropical Research Institute	F

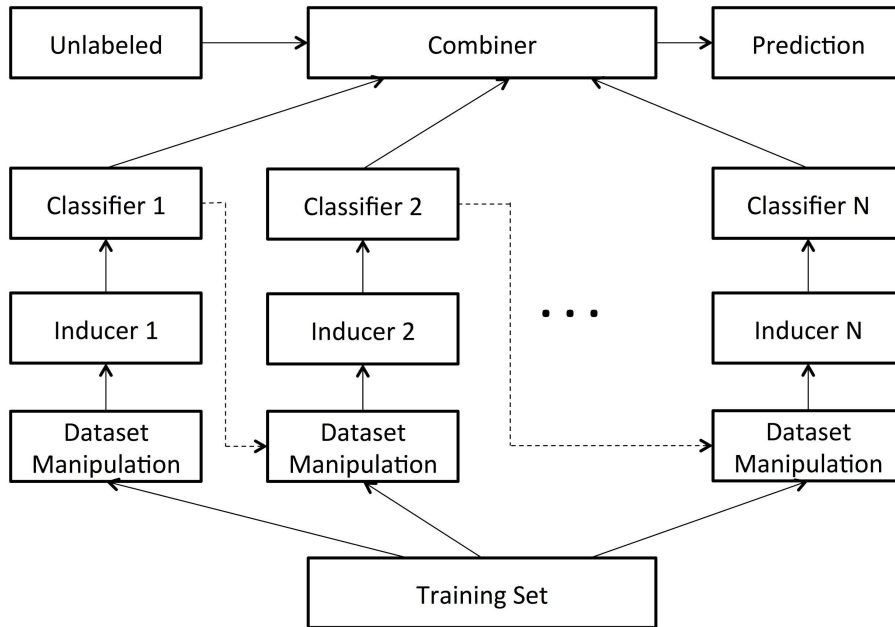
# Slot Filler Validation - Classification

- Slot Filler Validation is a binary classification task
  - Given a set of queries consisting of tuples of the form <entity, slot>
  - And a set of Slot Fillers for each query
  - Determine if a slot filler is True or False

# Slot Filler Validation - Classification

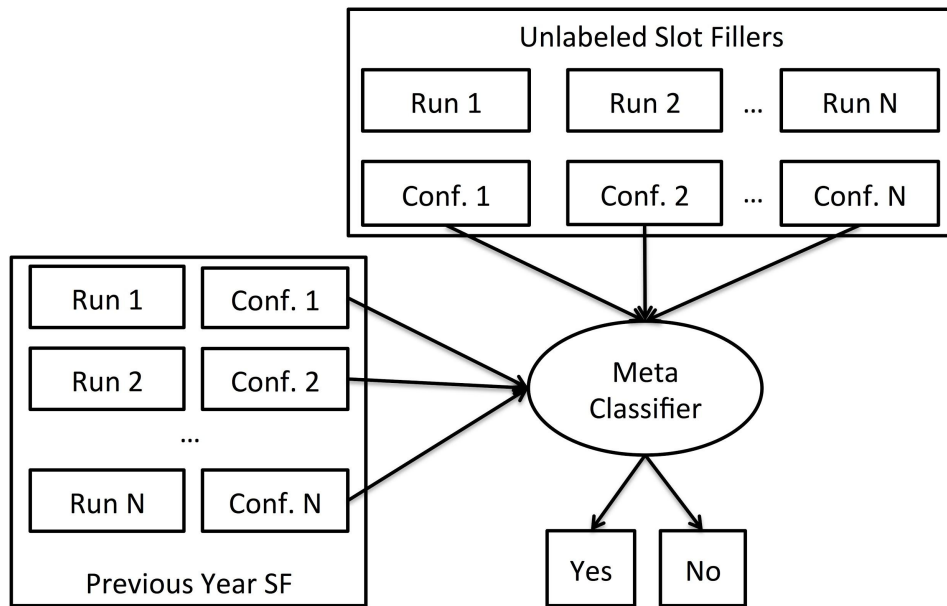
- Slot Filler Validation is a binary classification task
  - Given a set of queries consisting of tuples of the form <entity, slot>
  - And a set of Slot Fillers for each query
  - Determine if a slot filler is True or False
- A CSSF output is the output of such classifier
  - Ideal for ensemble classification
  - Aggregate the output of multiple classifiers
  - Outperform the original ones

# Ensemble Classification



- Ensemble methods have two main parts
  - **Inducer**: Selects the training data for each individual classifier
  - **Combiner**: takes the output of each classifier and combine them to formulate a final prediction

# Stacked Ensemble



Meta-level classifier that takes the output of other models as input and estimate their weights



# Stacked Ensemble

- Requires labeled data
  - Available from 2013 and 2014 SF and SFV
- Training Strategy
  - Learn from previous year performance
  - 2013-2014: 7 teams
  - 2014: 12 teams

# Stacked Ensemble

- Requires labeled data
  - Available from 2013 and 2014 SF and SFV
- Training Strategy
  - Learn from previous year performance
  - 2013-2014: 7 teams
  - 2014: 12 teams
- All runs that can not be fit into the classifier are discarded!
  - Leave out extra evidence
  - ... From potentially well ranked systems

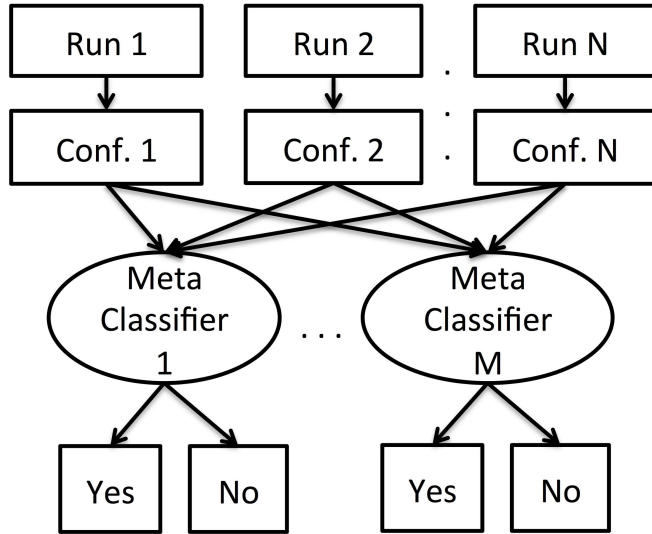
# Stacked Ensemble - not enough!

Rank	TEAM ID	0-HOP F1	1-HOP F1	ALL F1
9	SFV2015_SF_03_1	0.3457	0.1154	0.2718
14	SFV2015_KB_16_2	0.2633	0.1655	0.2247
16	SFV2015_SF_18_1	0.292	0.0972	0.2245
24	SFV2015_SF_08_4	0.2669	0.0976	0.2102
31	SFV2015_SF_02_1	0.1883	0.1299	0.1649
34	SFV2015_SF_06_1	0.2351	0	0.1595

Rank	TEAM ID	0-HOP F1	1-HOP F1	ALL F1
39	SFV2015_KB_10_1	0.1834	0.0952	0.1474
45	SFV2015_KB_09_1	0.0965	0.0791	0.0899
47	SFV2015_SF_13_2	0.1225	0	0.0892
56	SFV2015_SF_07_1	0.0512	0	0.0353
63	SFV2015_KB_11_1	0.019	0	0.0121
64	SFV2015_SF_17_1	0.019	0	0.0121

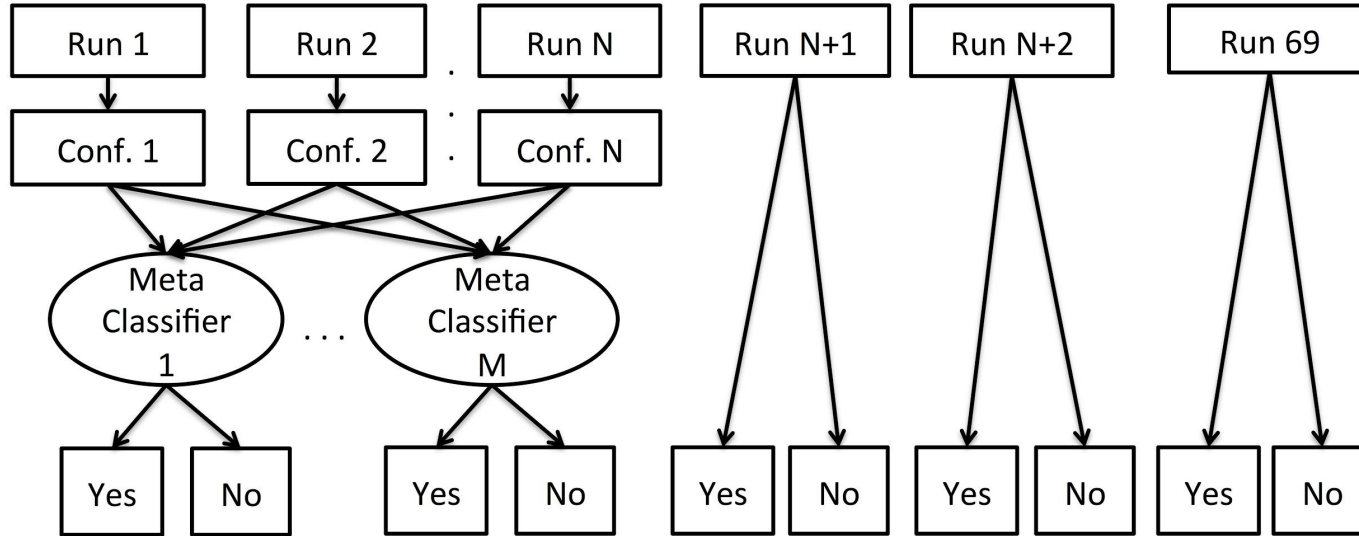
F1 score ranking of 2014-2015 teams.

# Consensus Maximization Fusion



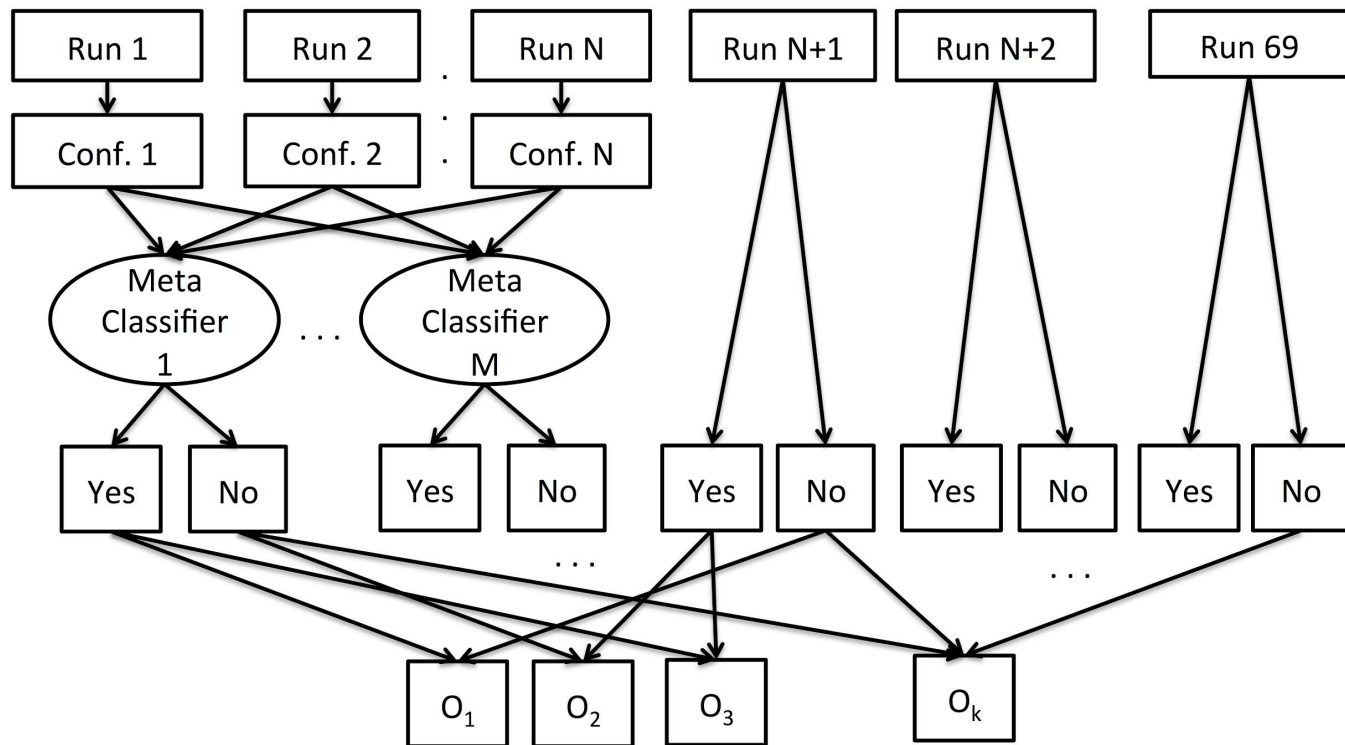
Augment stacked ensemble model by adding more meta-classifiers

# Consensus Maximization Fusion



Add runs that can not fit into the stacked ensemble method. We treat these runs as 2-Class Clusters

# Consensus Maximization Fusion



# Consensus Max. Fusion - Example

- Consider the following queries
  - O1 = (Marion Hammer, per:title, president)
  - O2 = (Dublin, gpe:headquarters\_in\_city,trinity college)




[More Images](#)

## Marion Hammer

Marion P. Hammer was the first female President of the National Rifle Association.  
[Wikipedia](#)

**Born:** [Columbia, SC](#)



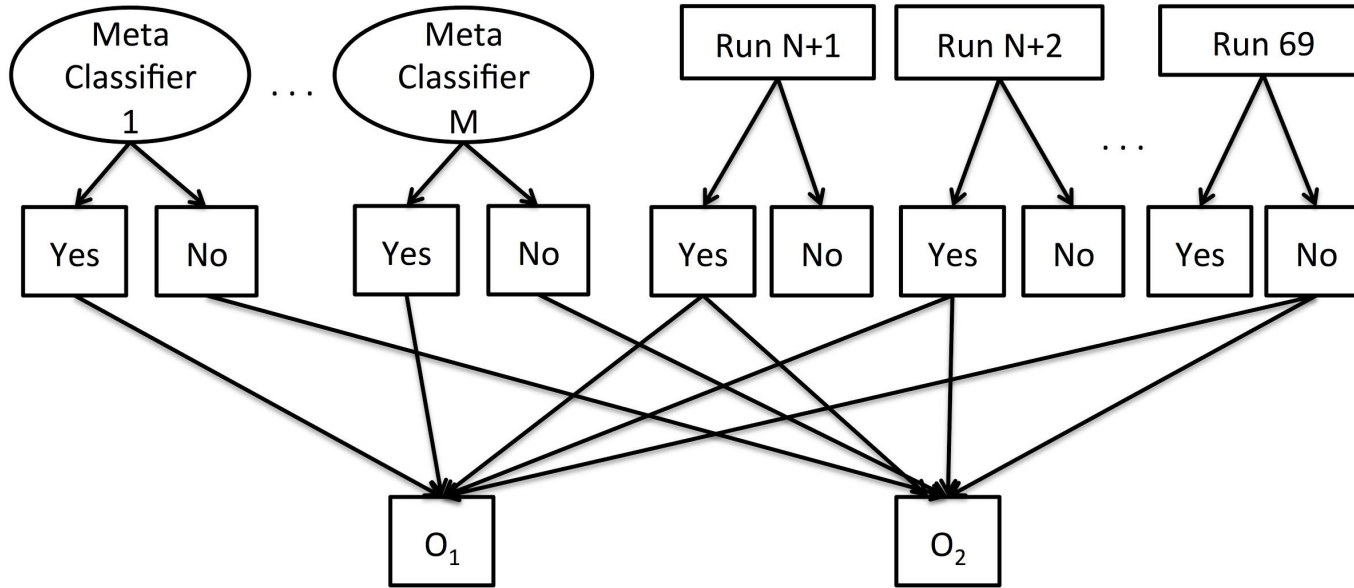
## Trinity College, Dublin ★

College in  
Dublin,  
Republic of  
Ireland

[Website](#) [Directions](#)

Trinity College, known in full as the College of the Holy and Undivided Trinity of Queen Elizabeth near Dublin, is a research university and the sole

# Consensus Max. Fusion - Example



Meta-Classifiers: 6 Yes – 0 No  
Clusters: 46 Yes - 16No

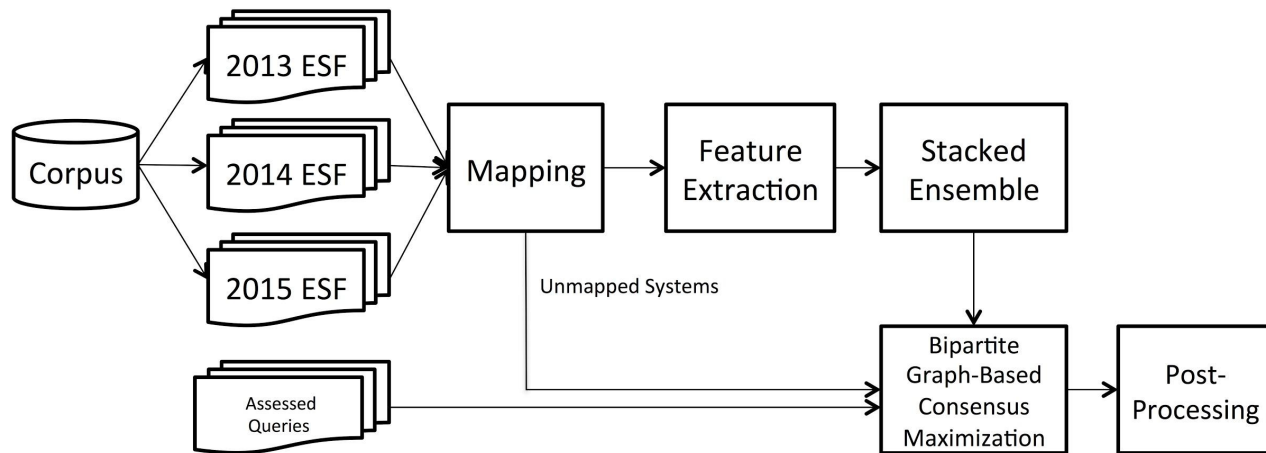
Meta-Classifiers: 0 Yes – 6 No  
Clusters: 34 Yes - No 28



# Consensus Max. Fusion

- Combine outputs of multiple supervised and unsupervised models for better classification.
- The predicted labels should agree with the base supervised models but adds unsupervised evidence.
- Model combination at output level is needed in KBP applications where there is no access to individual extractors.

# Consensus Maximization Fusion Pipeline



# Mapping

- Runs from teams that participated in previous years are mapped together and ranked using the corresponding assessments.
- 2015 runs, are ranked based on the small assessment file provided for the task.
- The best run of each mapped team is then passes to the feature extraction module.
- All other runs are passed directly to BGCM.

# Feature Extraction

- Same as the SFV Stack Ensemble System
  - Probabilities
  - Relation
  - Provenance

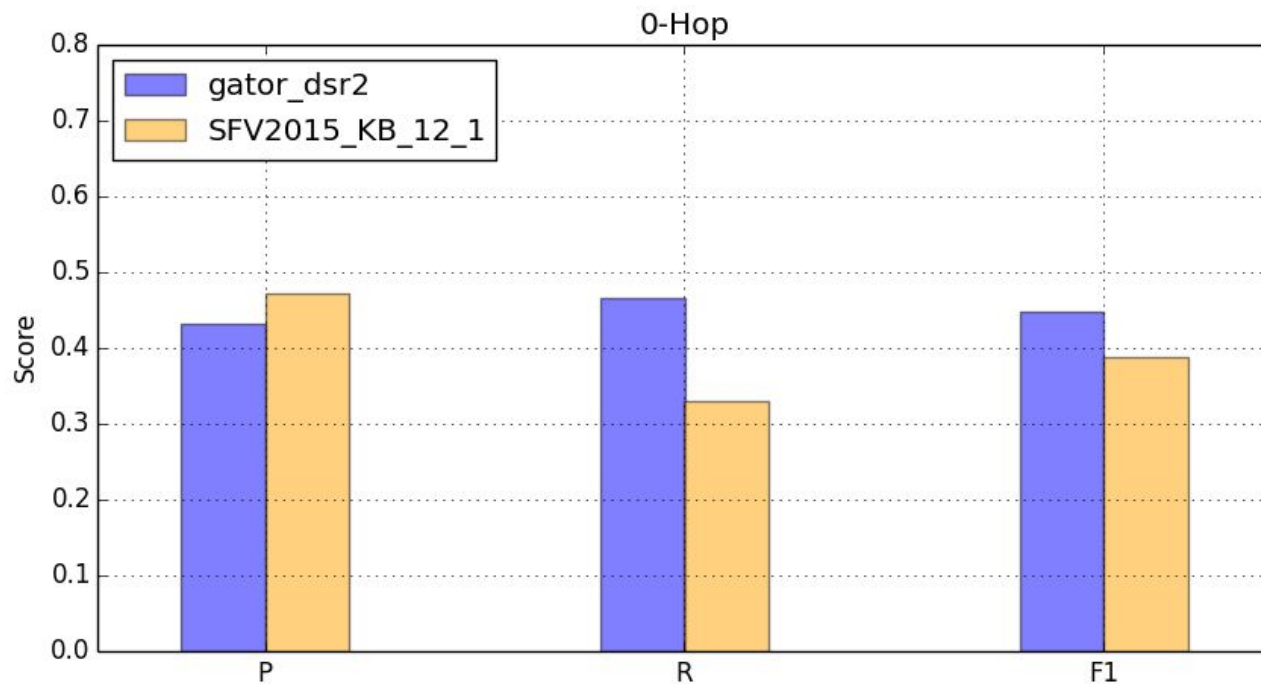
# Post-processing

- Filter ensemble of all 0-hop queries
  - Enforce single-values relations by selecting the one with highest probability
  - For every slot filler classified as true, select the provenance of the slot filler with highest probability.
- For every 1-hop query in the ensemble
  - Enforce its 0-hop result is in the ensemble

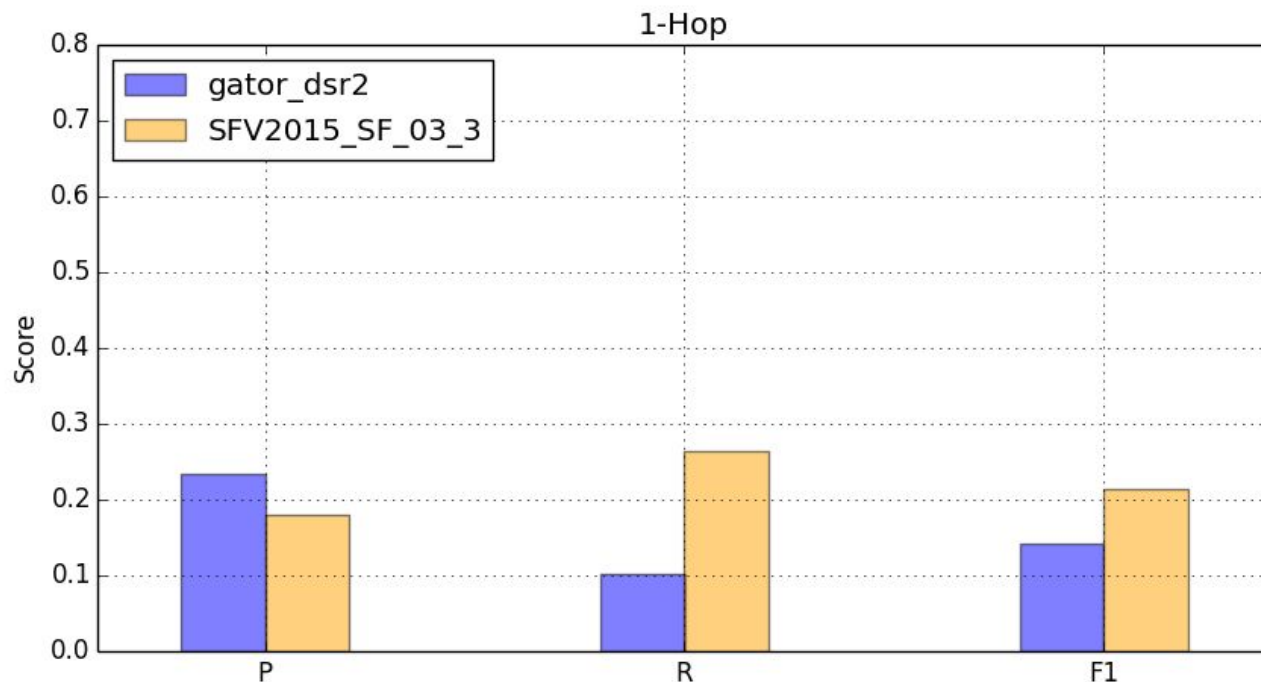
# Submitted Runs

- 2013-2014: Run 1
  - Meta-classifiers trained with samples from 7 teams.
  - BGCM: 6 meta-classifiers and 62 runs
- 2014: Run 2
  - Meta-classifiers trained with samples from 12 teams.
  - BGCM: 6 meta-classifiers and 57 runs
- Run 3
  - Use all meta classifiers from Runs 1 and 2
  - BGCM: 12 meta-classifiers and 57 runs

# Results - 2015 CSSF

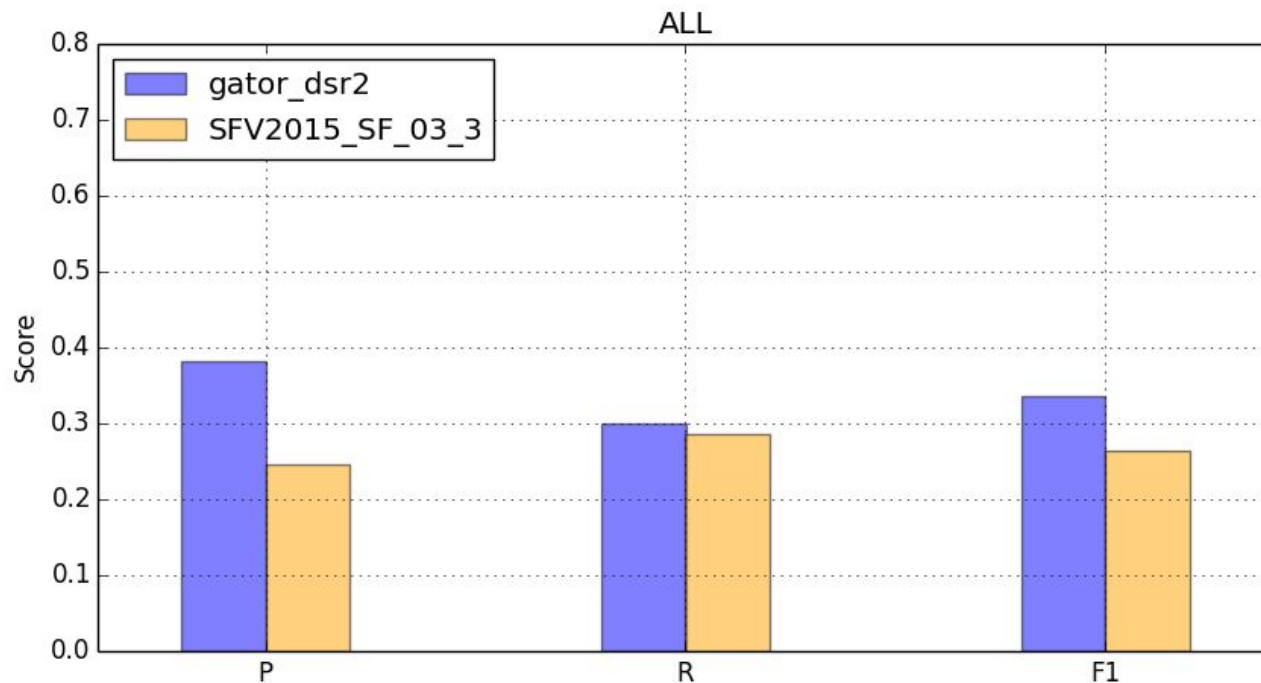


# Results - 2015 CSSF

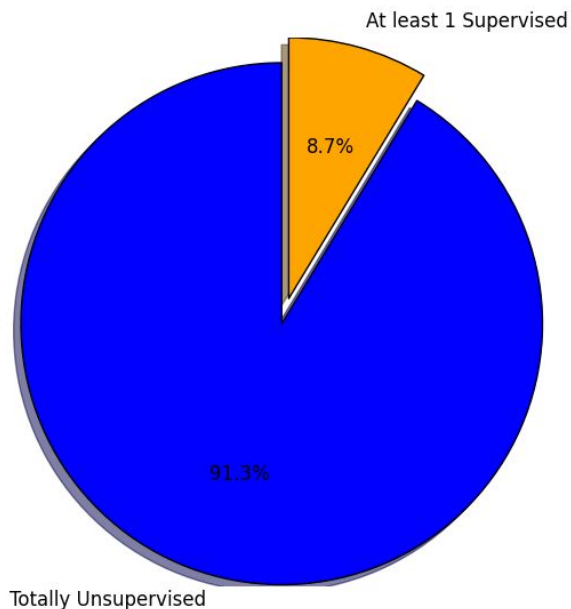




# Results - 2015 CSSF

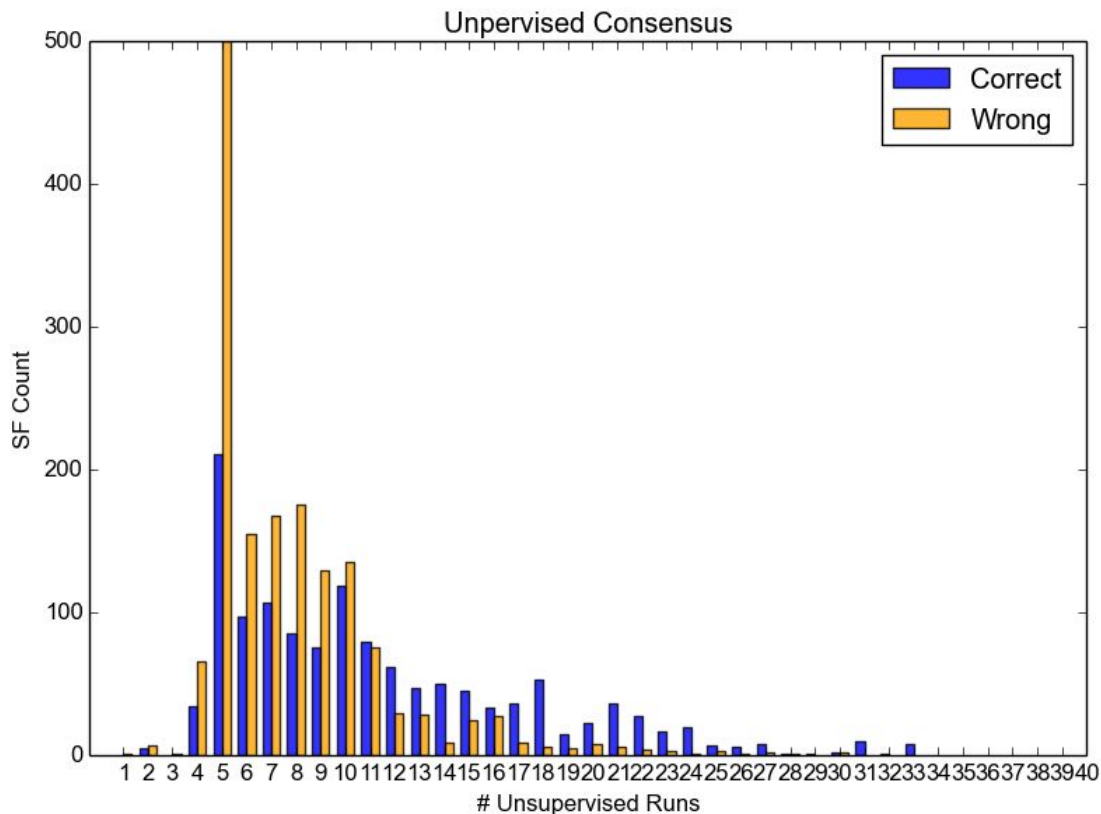


# Analysis Run 2



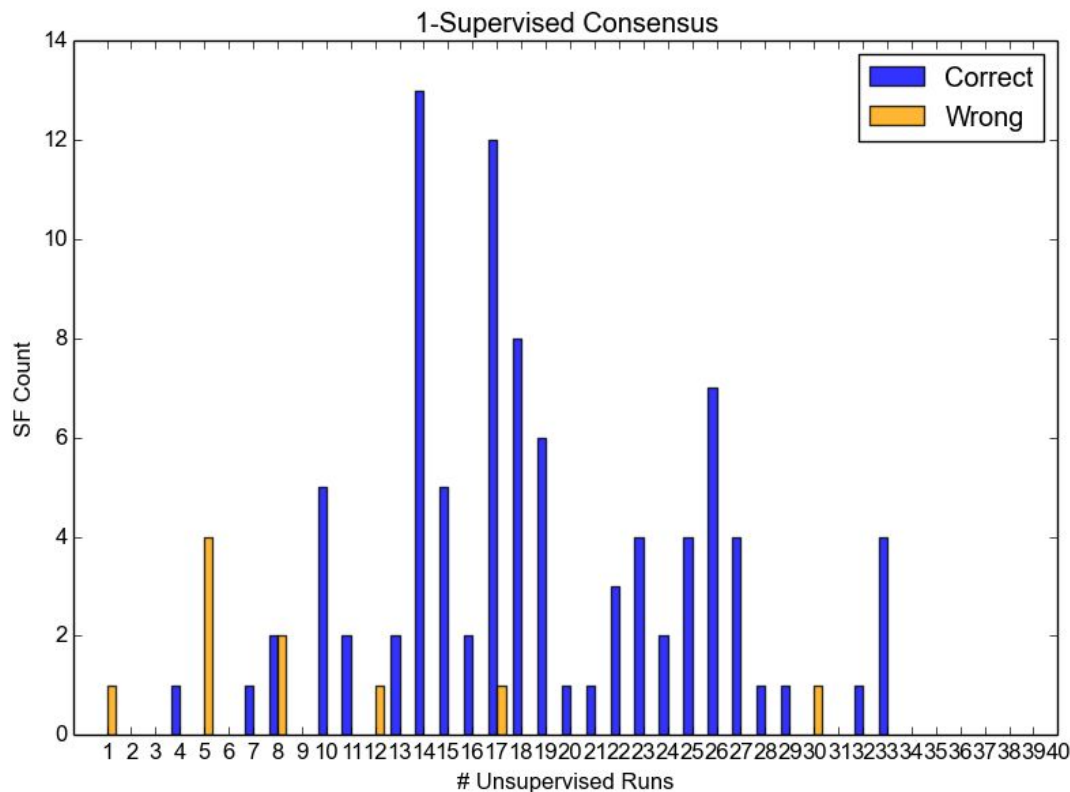
The majority of the slot fillers included in our best run come from unsupervised consensus

# Analysis Run 2



- Answers come from unsupervised consensus
  - All supervised outputs classified them as negative
  - Not enough evidence
- As more unsupervised runs reach consensus, there are more correct than incorrect fillers.
- The Recall of the system is improved

# Analysis Run 2



- At least one stacked ensemble model classified as positive.
- Supervised evidence helps improve precision.
- The higher the consensus with the unsupervised clusters the system filters better.

**Questions?**