





**ISTANBUL TECHNICAL UNIVERSITY ★ INFORMATICS INSTITUTE**

**UNSUPERVISED ACTIVE LEARNING FOR VIDEO ANNOTATION**

**M.Sc. THESIS**

**Emre DEMİR**

**Informatics Institute**

**Computer Sciences Department**

**MAY 2015**



**ISTANBUL TECHNICAL UNIVERSITY ★ INFORMATICS INSTITUTE**

**UNSUPERVISED ACTIVE LEARNING FOR VIDEO ANNOTATION**

**M.Sc. THESIS**

**Emre DEMİR  
(704101003)**

**Informatics Institute**

**Computer Sciences Department**

**Thesis Advisor: Prof. Zehra ÇATALTEPE**

**MAY 2015**



**İSTANBUL TEKNİK ÜNİVERSİTESİ ★ BİLİŞİM ENSTİTÜSÜ**

**VIDEO ETİKETLEME İÇİN AKTİF ÖĞRENME**

**YÜKSEK LİSANS TEZİ**

**Emre DEMİR  
(704101003)**

**Bilişim Enstitüsü**

**Bilgisayar Bilimleri Bölümü**

**Tez Danışmanı: Prof. Zehra ÇATALTEPE**

**MAYIS 2015**





**Emre DEMİR**, a M.Sc. student of ITU Informatics Institute 704101003 successfully defended the thesis entitled “**UNSUPERVISED ACTIVE LEARNING FOR VIDEO ANNOTATION**”, which he/she prepared after fulfilling the requirements specified in the associated legislations, before the jury whose signatures are below.

**Thesis Advisor :**    **Prof. Zehra ÇATALTEPE**    .....

Istanbul Technical University

**Jury Members :**    **Assist. Prof. Yusuf YASLAN**    .....

Istanbul Technical University

**Assist. Prof. Arzucan ÖZGÜR**    .....

Boğaziçi University

**Date of Submission :**    **May 2015**

**Date of Defense :**    **28 May 2015**



*To my wife Melis,*



## **FOREWORD**

I would like to thank Prof.Dr. Zehra ataltepe for her assistance, motivation and involvement in every step throughout the thesis.

I would like to thank Umit Ekmekci and Mateusz Budnik for their help, ideas and collaboration.

I thank to my beloved wife Melis zgür etinkaya Demir for her encouragement and support.

Work in this paper is partially supported by the Scientific and Technological Research Council of Turkey (TUBITAK) project 112E176.

May 2015

Emre DEMİR



## TABLE OF CONTENTS

	<u>Page</u>
<b>FOREWORD</b> .....	<b>ix</b>
<b>TABLE OF CONTENTS</b> .....	<b>xi</b>
<b>ABBREVIATIONS</b> .....	<b>xiii</b>
<b>LIST OF TABLES</b> .....	<b>xv</b>
<b>LIST OF FIGURES</b> .....	<b>xvii</b>
<b>SUMMARY</b> .....	<b>xix</b>
<b>ÖZET</b> .....	<b>xxi</b>
<b>1. INTRODUCTION</b> .....	<b>1</b>
1.1 Thesis Organization.....	2
<b>2. NOTATION AND DATASET</b> .....	<b>3</b>
2.1 Notation .....	3
2.2 Dataset .....	3
<b>3. RELATED WORK</b> .....	<b>7</b>
3.1 Active Learning .....	7
3.1.1 Active Learning Scenarios.....	8
3.1.1.1 Query Synthesis .....	9
3.1.1.2 Stream Based Selective Sampling .....	9
3.1.1.3 Pool-Based Sampling.....	10
3.1.2 Sampling Strategies .....	10
3.1.2.1 Uncertainty Sampling .....	10
Least Confident .....	11
Margin .....	11
Entropy .....	11
3.1.2.2 Query by Disagreement .....	12
3.1.2.3 Query by Committee.....	12
3.2 Cluster Matching Metrics .....	13
<b>4. METHODOLOGY</b> .....	<b>19</b>
4.1 Clustering Algorithms .....	19
4.1.1 Agglomerative Clustering.....	21
4.1.2 K-Medoid Clustering.....	21
4.2 Query Instance Selection Strategies .....	23
4.2.1 Cluster Selection.....	25
4.2.1.1 Big Cluster First.....	25
4.2.1.2 The Most Disagree Selection.....	25
4.2.1.3 Hybrid Cluster Selection.....	26

4.2.2 Instance Selecting Strategy.....	27
<b>5. EXPERIMENTS AND RESULTS .....</b>	<b>29</b>
5.1 Evaluation Metrics.....	29
5.2 Results .....	30
<b>6. CONCLUSIONS AND RECOMMENDATIONS .....</b>	<b>37</b>
<b>REFERENCES.....</b>	<b>39</b>
<b>CURRICULUM VITAE .....</b>	<b>44</b>



## ABBREVIATIONS

<b>BCPS</b>	: Big Cluster Pair First Selection
<b>BCS</b>	: Big Cluster First Selection
<b>HHC</b>	: Hard Hybrid Certainty Selection
<b>HHU</b>	: Hard Hybrid Uncertainty Selection
<b>MDS</b>	: Most Disagree Selection
<b>DS</b>	: Disagreement Score
<b>DSL</b>	: Disagreement Score List
<b>HYB</b>	: Soft Hybrid Selection
<b>H-HYB</b>	: Hard Hybrid Selection
<b>KL</b>	: Kullback-Leiber Divergence
<b>OCR</b>	: Optical Character Recognition
<b>QDB</b>	: Query By Disagreement
<b>QDC</b>	: Query By Committee
<b>SHS</b>	: Soft Hybrid Score
<b>VI</b>	: Variation of Information
<b>ST</b>	: Speaker Track
<b>FT</b>	: Face Track
<b>ON</b>	: Overlaid Name
<b>SS</b>	: Speaker Score with Speaker Annotation
<b>FS</b>	: Face Score with Speaker Annotation
<b>FF</b>	: Face Score with Face Annotation
<b>SF</b>	: Speaker Score with Face Annotation



## LIST OF TABLES

	<u>Page</u>
<b>Table 2.1</b> : The video list and the size of distance matrices which are inputs for our algorithms. ....	6
<b>Table 4.1</b> : A sample cluster pair similarity matrix for agglomerative and K-Medoid algorithms.....	22



## LIST OF FIGURES

	<u>Page</u>
<b>Figure 2.1</b> : Example frames for face segmentation and OCR [1]......	4
<b>Figure 3.1</b> : Detecting safe shaped food threshold $\theta$ with binary search.....	8
<b>Figure 4.1</b> : The architecture of the system. ....	20
<b>Figure 4.2</b> : The illustration of Agglomerative Clustering.....	21
<b>Figure 4.3</b> : Pseudo-code of Agglomerative Clustering .....	21
<b>Figure 4.4</b> : Pseudo-code of K-Medoid algorithm. ....	22
<b>Figure 4.5</b> : The stable matching algorithm. ....	24
<b>Figure 5.1</b> : Face score with face annotation.....	31
<b>Figure 5.2</b> : Speaker score with face annotation. ....	31
<b>Figure 5.3</b> : Face score with with speaker annotation. ....	31
<b>Figure 5.4</b> : Speaker score with with speaker annotation.....	31
<b>Figure 5.5</b> : Face score with face annotation.....	33
<b>Figure 5.6</b> : Speaker score with face annotation. ....	33
<b>Figure 5.7</b> : Face score with with speaker annotation. ....	33
<b>Figure 5.8</b> : Speaker score with with speaker annotation.....	33
<b>Figure 5.9</b> : Face score with face annotation (lower is better). ....	34
<b>Figure 5.10</b> : Speaker score with face annotation (lower is better).....	34
<b>Figure 5.11</b> : Face score with with speaker annotation (lower is better).....	34
<b>Figure 5.12</b> : Speaker score with with speaker annotation (lower is better). ....	34



# UNSUPERVISED ACTIVE LEARNING FOR VIDEO ANNOTATION

## SUMMARY

When annotating complex multimedia data like videos, a human expert usually annotates them manually. Even though manual annotation achieves accurate results, it is a labor-intensive and time-consuming process. On the other hand, computational methods can annotate mass video data for indexing and searching with any or almost no help from human experts effortlessly and faster but they are probably more error prone solutions. The tradeoff between the costs in terms of labor, time and accuracy reveals Active Learning as a natural outcome. Active learning is one of the semi-supervised machine learning methods that benefits from the strongest properties of both manual and computational methods. In an active learning cycle, a learner algorithm discovers the underlying patterns in data and queries the human experts interactively for some informative decision points. It is used when labeled instances are insufficient and acquiring new labels is expensive or especially when unlabeled instances are abundant. In this study, we introduce an unsupervised active learning cycle structure in a flow, which includes clustering, stable matching between the created clusters, various unsupervised selection strategies for selecting the most uncertain and the most certain instances and querying the human annotators. We propose two new cluster selection methods, namely Most Disagreement Selection (MDS) and Hybrid Set Selection (HS), which is a hybrid of MDS and Big Cluster First [2] methods. For MDS and HS, we adopt the "Stable Marriage Problem" solution, in which a stable marriage problem is transformed into a cluster matching problem. We work on REPERE [1] video dataset, which is created for the problem of person identification in videos. Our study aims to identify who is speaking and who is on screen by using multi-modal data. We have evaluated the performance of selection strategies over active learning cycles using multimodality on 28 videos from 7 different TV programs. Each video has three different similarity matrices namely face-to-face, speech-to-speech and face-to-speech. We have run four experiments with regard to matrices in this order: face score for face track annotation, face score for speaker track annotation, speaker score for speaker track annotation and speaker score for speaker annotation.





## VIDEO ETİKETLEME İÇİN AKTİF ÖĞRENME

### ÖZET

Günümüzde dijital video ve fotoğraf içeriklerinin kolay üretilebilmesi ve bu üretilen videoların internet üstünden medya, sosyal medya ve video paylaşım siteleri gibi kanallar üstünden internet sitelerine son kullanıcıların da çok kolay yükleme yapabilmesi nedeniyle her geçen gün internet ortamında muazzam derecede içerik oluşmasına neden olmaktadır. Video verilerinin sürekli artması ve bu verilere kullanıcıların çok kolay erişmesi, bilgiye ulaşmayı oldukça arttırmaktadır ancak video verisinin artması ve videoların içeriklerinin etiketlenmesinin zor olmasından dolayı aynı zamanda doğru videolara erişim, istenilen videoyu arama, istenilen videoların arşivlenmesi vb. işlemleri zor kılmaktadır. Bunun nedenlerinden başında video verisini üreten ve/veya video verisini kayıt eden kullanıcıların video içeriklerini etiketlememesi veya eksik etiketlemesidir. Bu durum etiketlenmemiş videoların aranıp bulunmasını oldukça zorlaştırmaktadır. Ayrıca etiketlenmemiş, eksik etiketlenmiş veya yanlış etiketlenmiş video verileri aynı zamanda arama veya arşivleme gibi işlemlerde hatalı ve gürültülü sonuçlara neden olmaktadır.

Kullanıcıların videoları etiketlememesi veya düzgün etiketlememesinin nedenlerinden birisi videoların içeriklerinin etiketlenmesi işleminin zor olmasıdır. Tüm videolar baştan sona izlenip içeriklerinin ayrıntılı olarak etiketlenmesi gerekmektedir. Bu etiketler video içeriğinde hangi insanların olduğu gibi etiketler olabilirken, videoda kayıt edilen genel bir olay(basketbol, yangın videosu) veya spesifik bir olay da olabilir, dünya kupası finali gibi. Kullanıcıların etiketlemedikleri veya eksik etiketleme yaptıkları durumlarda video verisini etiketlenebilmesi için otomatik video etiketleme konusunda araştırma ve geliştirme yapılmaktadır.

Bu çalışmada otomatik video etiketleme yapacak makine öğrenmesi algoritmalarının eğitimi sırasında gereksinim duyduğu etiketlenmiş video verisinin azaltılmasına odaklanılmıştır. Denetimli makine öğrenmesi temelli otomatik video etiketleme sistemleri, öğrenme safası öncesinde insan uzmanlar tarafından etiketlenmiş belli sayıda örnek kullanır. Denetimli makine öğrenme algoritmalarının eğitim sürecince kullanacağı örnekler genellikle rastgele seçilir ve insan uzman tarafından etiketleme işlemi yapılır. Otomatik etiketleme sistemlerinin öğrenme safası için gereksinim duydukları bu etiketlenmiş örneklerin bir insan tarafından etiketlenmesi süreci çok uzun sürmektedir. çoğu zaman etiketlenecek video kare kare incelenip her kare etiketlenir. Etiketleme işleminin uzun sürmesi nedeniyle makine öğrenmesi algoritmalarının öğrenim sürecinde daha az etiketlenmiş örneğe gereksinim duyan yöntemler araştırılmakta ve geliştirilmektedir. Bu yöntemlerden bir tanesi yarı denetimli bir yöntem olan aktif öğrenmedir.

Denetimli öğrenme yöntemlerinde, hipotezin gereksinim duyduğu kadar örnek bir uzman tarafından etiketlenir ve bu örnekler tüm örnek havuzundan rastgele seçilir. Buna karşılık yarı denetimli yöntem olan aktif öğrenme algoritmaları çok az sayıda rastgele etiketlenmiş veri ile öğrenim sürecine başlar ve denetimli öğrenmeye göre daha az sayıda etiketlenmiş örnek ile denetimli makine öğrenmesi yöntemiyle aynı başarı seviyesine gelmeyi amaçlar. Bunun başarılabilmesi örnekleri rasgele seçmek yerine, sınıflandırıcı için en çok bilgi taşıyan örneklerin bulunmasına bağlıdır. En çok bilgi taşıyan örnekler teorik olarak karar sınırlarının çevresinde bulunur. İkili bir sınıflandırma problemi için anlatacak olursak, karar sınırlarının üstünde yer alan bir örneğin her iki sınıfa da ait olma olasılığı 0.5 dir. Olasılık 0.5 olduğu durumda entropi en yüksek değerini alır ve bu noktaların bulunması etiketlenecek örnek sayısında düşüşe neden olması beklenir. Bu noktada aktif öğrenme çözülmeye çalışılan problem için en çok bilgi taşıyan örnekleri bulmaya çalışır, teoride karar sınırlarının etrafında yer alan örnekler en çok bilgiye sahip olması beklense de önceki çalışmalar göstermiştir ki bazı hipotezler de karar sınırının en uzağında yer alan örnekler daha çok bilgi taşıyabilmektedir.

Biz bu çalışmada video etiketleme için aktif öğrenme yöntemini denetimsiz öğrenme problemi üstünde uyguladık. Böylece otomatik video etiketleme sisteminin rastgele etiketleme ile aynı seviyede başarıyı daha az sayıda etiketlenmiş video ile elde edilmesi amaçlandı.

Aktif öğrenme fikri 1980'li yıllarda çıkmış olup, uygulamaları günümüze kadar yaygınlaşmamıştır. Günümüzde aktif öğrenme ile ilgili çalışmalar hızlanmıştır ve bu hızlanmanın ana nedeni büyük veriler üstünde öğrenme algoritmalarının öğrenme safasını hızlandırmaktır. Bu çalışmada önerdiğimiz yöntem denetimsiz bir makine öğrenmesi olan kümeleme üstünde çalışmaktadır. çalışmanın detaylarında anlatıldığı gibi aktif öğrenme algoritmaları genellikle denetimli yöntemler için geliştirilmiş ve denetimli yöntemler üstünde uygulamaları daha yaygındır. Bu çalışma sırasında denetimsiz öğrenme için aktif öğrenme tekniği geliştirilerek problemleri çözülmüştür.

Önerilen yöntem denetimli aktif öğrenme yöntemi olan "anlaşmazlığa göre sorgulama" (query by disagreement) ve "belirsizi sorgulama" (uncertainty sampling) yöntemlerinin denetimsiz öğrenme için geliştirilmişidir.

Anlaşmazlığa göre sorgulama denetimli yöntemlerde uygulaması şu şekilde tanımlanır; aynı veri kümesi kullanılarak iki farklı sınıflandırıcı eğitilir, bu iki farklı sınıflandırıcı aynı tip sınıflandırıcının fazla eğitilmiş bir hipotezi ve az eğitilmiş başka bir hipotezi olabileceği gibi iki farklı tip sınıflandırıcı da olabilir. İki farklı sınıflandırıcı eğitildikten sonra bu iki sınıflandırıcı etiketlenmemiş örnekleri ayrı ayrı etiketler ve yani her bir örneğin farklı sınıflandırıcılardan gelen iki adet etiketi olur. Bir örneğe atanan iki farklı etiket birbirinden farklı ise yani sınıflandırıcılar anlaşmazlığa düşmüş durumdaysa, bu örnek etiketlenmek üzere insan uzamana sorulur. Böylece örneklerin dağılımında sınıflar arasındaki karar sınırına en yakın örneklerin bulunması amaçlanır çünkü iki sınıflandırıcının en çok anlaşmazlığa düştüğü noktalar karar sınırına yakın olanların olması beklenir.

Bir başka aktif öğrenme yöntemi en basit ve en çok kullanılan denetimli yöntem olan 'belirsizliği sorgulama'. Bu yöntemde karar sınırına en yakın örnekler seçilir ve etiketlenmesi için insan uzmana sorulur. Biz çalışmamızda anlatılan bu iki aktif

öğrenme yöntemini de denetimsiz öğrenme için geliştirip REPERE verisi üstünde uyguladık.

Anlaşmazlığa göre sorgulama yönteminin denetimsiz öğrenme probleminde geliştirilmesinin ve problemlerinin çözümünün ilk adımı olarak iki adet denetimsiz kümeleme yöntemi eğitilmiştir. Kümeleme yöntemi olarak yığılmal kümeleme (agglomerative clustering) ve K-medoid kümeleme yöntemleri kullanılmıştır.

Denetimli yöntemde anlaşmazlığa göre sorgulamada sınıflandırılar test aşamasında etiket atadığından anlaşmazlık direkt anlaşılıp insan uzamana sorulurken, denetimsiz yöntemde etiketler bulunmamakta, etiketsiz kümeler bulunmakta. Bu durumda iki kümeleme yönteminin kümelerini karşılaştırabilmek için aynı etiketi temsil eden kümelerin birbiriyle işleşmesi gerektiği düşünülmüş ve eşleşme problem çözülmüştür. Kümelerin eşleştirme işlemi için yığılmal kümeleme yönteminden gelen her bir küme için K-medoid yönteminin çıktısı olan her bir küme ile ortak eleman sayısı bulunup toplam eleman sayısına bölünmüş ve kümeler arasında bir anlaşmazlık değeri matrisi oluşturulmuştur. Kümeler birlerine göre değerlendirildikten sonra oluşturulan benzerlik matrisi kullanılarak birbirine en uygun kümeler kararlı evlilik probleminin (stable marriage problem) çözüm algoritması uygulanmış ve en uygun küme çiftlerinin birbirine atanması garanti altına alınmıştır. İki kümeleme yönteminin ürettiği kümeler birbirlerine atanma işleminden sonra aralarında en fazla anlaşmazlık olan küme çifti içlerinden sorgu yapılması amaçlı seçilmiştir. Biz bu yöntemi "en çok anlaşmazlığı sorgulama" (MDS) olarak adlandırdık. Bu aşamada hangi küme çifti içinden örnek seçileceği belirlendi ancak kümeler içinde hangi örneğin sorulacağını seçmedik.

Seçilen küme çifti içinden örnek sorgulanması için örneğin seçiminde belirsizliğe -uncertainty sampling- göre sorgulama yöntemini küme için uyguladık. Bu adımda seçilen küme çifti içinde ki örnekler tek tek taranarak diğer örneklere toplam uzaklığı en fazla olan örneği seçtik. Diğer örneklere toplam uzaklığın en fazla olduğu örnek bize küme içinde kümenin sınırına en yakın olan yani belirsizliği en fazla olan örneği verir. Deneyler sırasında görüldü ki küme içinde en belirsiz noktaları, yani olasılığı en düşük noktaları seçmenin yanı sıra olasılığı en yüksek olan noktaları yani kümenin merkezine yakın olan, olasılığı en yüksek olan noktaları seçip sorgulamak da başarıyı arttırmakta. Bu nedenle belirsizliği en düşük olan noktalar seçilip sorgulama yapıldığı deneylerde bu çalışmada yapıldı.

Bu çalışmada önerilen yöntem [2] çalışmasında önerilen "Big Cluster First" (BCS) yöntemiyle karıştırıldı ve iki yöntemin de güçlü olduğu bölümlerin olduğu görüldü. Bu nedenle BCS ve MDS yöntemleri birleştirilerek melez yöntemde önerilmiştir.

Deneylerde aktif öğrenme algoritmaları REPERE video verisi üstünde 4 farklı sınıflandırma problem üstünde yapılmıştır. Burada amaç videoda o anda kimin konuştuğunun düzgün etiketlenmesidir. Sınıflandırma problemlerinden ilki konuşmacı yani ses verisinden gelen öz nitelikler kullanılarak, videoda ki konuşmacının etiketlenmesi. İkincisi problemde ses verisinden gelen öz nitelikler kullanılarak videoda gözüken kişinin etiketlenmesi. üçüncü problemde görüntü verilerinden gelen öz nitelikler kullanılarak videoda gözüken kişinin etiketlenmesi ve son sınıflandırma probleminde görüntü verisinden gelen öz nitelikler kullanılarak o anda ki ses verisinin

etiketlenmesi. Anlaşılacağı üzere önerilen aktif öğrenme yöntemleri çok kipli problem üstünde uygulanmıştır.

Deney sonuçları 28 video için mikro F-measure, makro F-measure hesaplanarak değerlendirilmiştir. Ayrıca yöntemlerin başarıları kendi aralarında derecelendirilmiştir ve en iyi yöntem belirlenmeye çalışılmıştır. Deney sonucunda görülmüştür ki en başarılı yöntem ufak farkla da olsa kesinliği en yüksek örneklerin sorgulandığı melez yöntem olduğudur.

## 1. INTRODUCTION

In the last decade, video recording technologies have moved from analog to digital and video recording devices have proliferated among wide range of users. Another evolution is proliferation of high speed internet. Many regular internet users started to create videos after digital video evolution and they can share these videos over the internet. In addition, video streaming service providers, such as Youtube, personal blogs, news web sites, also contribute new data every day. According to the Youtube statistics, 100 hours of video are uploaded to Youtube every minute and over 6 billion hours of video are watched on Youtube every month [3]. The enormous amount of video data causes challenges in search and retrieval operations. In order to index and search videos with the required contents effectively, annotation process is applied on them. There are two main video annotation techniques; manual annotation, in which a human expert annotate a video manually and automated video annotation, in which various computational methods annotate a video. In manual annotation approach, an expert watches whole videos frame by frame and annotates topics, occurred events, persons, types (e.g music, news etc) and other informative data. On the other hand, computational methods may annotate many videos automatically without any or less help from human. Automated video annotation methods are faster but more error prone than the manual method. Even though manual approach is the most accurate one in terms of video annotation in many cases, annotating immense quantity of videos manually is a labor-intensive and time-consuming process.

Automatic annotation techniques such as Active Learning offers various solutions to overcome the excessive cost of manual annotation. Active learning is a semi-supervised machine learning technique that aims to reduce annotation/labelling costs. Conventional supervised or unsupervised machine learning methods use random annotation for learning phase. In active learning, an algorithm annotates videos automatically but asks for the labels of the most informative instances to human expert

in order to learn data distribution more accurate. By this way, experts may train learning model with less annotation effort.

Various active learning methods have been proposed to video annotation problem. One of the studies [4] applies active learning for video annotation by comparing uncertainty sampling, the most probable sampling and random sampling to video indexing. A study [5] proposes video retrieval and annotation system called LIGVID which uses two active learning methods: 'relevance sampling' and 'uncertainty sampling'. Another study integrates SVM based active learning for feature selection to solve the text classification problem [6]. However, active learning for feature selection fails in that study, because of the use of a wrong feature reduction technique called GainRatio Feature Selection [6].

In a study [7], active learning is applied on networked data, of which nodes are 'papers' and links are 'references to other papers'. It uses a method based on query by disagreement and reduces paper annotation costs for classifying research papers. A study [8] extends the traditional active learning framework by including feedback on features alongside labeling the instances. It focuses on the effects of feature selection and human feedback for features in the setting of text categorization and applies uncertainty sampling based methods.

This study is a sub-project of the project CAMOMILE [9], which targets to produce an annotation framework for 3M data where the letter 'M' stands for multimodal, multimedia and multilingual. In this study, we propose a cluster based unsupervised active learning approach as a selection strategy on REPERE [1] video dataset, which is created for the person identification problem in videos. Our study aims to identify who is speaking and who is on screen.

## **1.1 Thesis Organization**

The rest of the thesis is organized as follows: Section 2 includes the notation for the study and a background of used dataset. Whereas, Section 3 includes related works on active learning and cluster matching, Section 4 introduces the methods. Section 5 includes the experiments and concludes the study.

## 2. NOTATION AND DATASET

### 2.1 Notation

Literally, the finite data set is  $X = x_1, x_2, \dots, x_n$  where the cardinality is  $|X| = n$ . The cluster set  $C = \{C_1, C_2, \dots, C_k\}$  represents the cluster sets of the set  $X$  with the assumption  $|C_i| > 0$  for all  $i = 1, 2, \dots, k$ . The set of all clustering of  $X$  is denoted as  $B(X)$ .  $C'$  is the second clustering for  $X$  where  $C' = \{C'_1, C'_2, \dots, C'_\ell\} \in S(X)$ .  $M = (m_{ij})$  is the confusion matrix of each clustering pair  $C, C'$ . The intersection between  $C_i$  and  $C'_j$  is a  $k \times \ell$  matrix where the  $i^{th}$  element gives the number of the elements in the intersection of  $C_i$  and  $C'_j$ .

$$m_{ij} = |C_i \cap C'_j|, 1 \leq i \leq k, 1 \leq j \leq \ell \quad (2.1)$$

The product of the clustering  $C_i$  and the clustering  $C'_j$  is denoted as  $C_i \times C'_j$  and defines the coarsest common refinement of those two clusterings.

$$C \times C' = \{C_i \cap C'_j | C_i \in C, C'_j \in C', C_i \cap C'_j \neq \emptyset\} \quad (2.2)$$

### 2.2 Dataset

We use video dataset from the REPERE challenge, which aims to provide a benefit to research on person identification in videos for multimodal conditions. The REPERE challenge tries to find answers to questions “Who is speaking?” “Who is present in the video?” and etc. by use of various information on speech and image modals. The REPERE Corpus occurs from 28 videos, which include 7 different types of shows such as news, talk show etc. and various numbers of participants from 3 persons to dozens in a video. In addition, the length of the videos has a range from 3 to 30 minutes,



**Figure 2.1:** Example frames for face segmentation and OCR [1].

which naturally causes various numbers of annotations for each video approximately 20 to 100 frames.

As distinct from other learning approaches, active learning interactively asks for annotation from an expert depending upon a selected strategy. So, an initial set of annotated data is needed for these queries. In general, embedded texts in videos can point the name of a represented speaker. However, it is expensive and time-consuming to obtain text annotations manually. Therefore, we use the extracted overlaid texts in videos by an Optical Character Recognition (OCR) system [10].

The used text detection method is adopted from the study [11]. It includes a Sobel filter and erosion/dilation operations that are employed for coarse detection following by filtering out false positive text boxes. Moreover, Google OCR system called Tesseract is used for text recognition. After the OCR system detects and extract texts automatically, human annotators verify and improve the results.

Speech and face segmentations processes are applied on videos to gather feature set for both training and testing. For Speech segmentation, speech tracks are gathered by splitting signals into acoustically homogeneous segments. For each video, the similarity matrix between speech tracks is calculated by using BIC criterion [12] with single full covariance Gaussians. The similarities are normalised into the interval [0,1].

For face segmentation, the first and the fifth frames of each shot is scanned to find a face of which initial pose can be frontal, half-profile or profile. Moreover, the segmentation approach uses online tracking, which determines the current position and the location



of a face in a frame by using the information of previous frame. After gathering a face image, 9-point mesh (an eye: 2 points, a nose: 3 points and lips: 2 points) is imposed on it. Then, a 490 dimensional feature is calculated with HOG descriptor. An average descriptor is calculated by using found features of a face on the same shot. For each face track, the number of feature dimension is decreased from 490 to 200 with the help of LDML approach [13]. At last, the similarity matrix between face tracks is calculated with Euclidean distance. The similarities are normalized into the interval [0,1].

The third similarity matrix occurs from correlation scores between the faces and the speakers in the interval [0,1] to build a multimodal clustering. Too small faces are eliminated for this distance matrix. The change in the color histograms of the region of a speaker's lips indicates person talking. Moreover, the size of a face and proximity to the center of the screen give us the clue that the current speech track is associated to the current face track.

The test set of REPERE challenge includes 1229 annotated frames from three hours videos. Instead of annotating all videos frame by frame, every 10 seconds a frame is selected and annotated by human annotators.

Literally, each instance of the dataset  $X = x_1, x_2, \dots, x_n$  represents a face track and  $Y = y_1, y_2, \dots, y_n$  represents a speech track. An instance  $x_i$  has 200 dimensional feature vector  $x_i = (x_i^1, x_i^2, \dots, x_i^{200})^T$ . Likewise an instance  $y_j$  has feature vector  $y_j = (y_j^1, y_j^2, \dots, y_j^m)^T$ . The distance matrices  $D_f$ ,  $D_s$  and  $D_{fs}$  represents the distances, respectively, between face-to-face, speech-to-speech and face-to-speech tracks in a video. The size of the distance matrices and speaker track number  $ST$ , face track number  $ST$  and overlaid name numbers  $ON$  are given in the Table 2.1.

**Table 2.1:** The video list and the size of distance matrices which are inputs for our algorithms.

<b>Video Name</b>	$D_f$	$D_s$	$D_{fs}$	$FT$	$ST$	$ON$
BFMTV BFMStory 2012-01-10	483 × 483	348 × 348	1020	759	372	36
BFMTV BFMStory 2012-01-23	507 × 507	329 × 329	1037	782	338	34
BFMTV BFMStory 2012-02-14	498 × 498	342 × 342	1024	753	348	33
BFMTV BFMStory 2012-02-20	566 × 566	330 × 330	1104	871	341	30
BFMTV CultureEtVous 2012-01-13	40 × 40	36 × 36	94	69	38	3
BFMTV CultureEtVous 2012-01-16	92 × 92	39 × 39	178	163	42	5
BFMTV CultureEtVous 2012-01-17	30 × 30	44 × 44	70	48	49	5
BFMTV CultureEtVous 2012-01-18	50 × 50	45 × 45	128	102	47	6
BFMTV CultureEtVous 2012-01-19	49 × 49	38 × 38	101	4	78	41
BFMTV CultureEtVous 2012-02-14	59 × 59	48 × 48	136	9	104	48
BFMTV CultureEtVous 2012-02-15	94 × 94	34 × 34	183	4	180	34
LCP CaVousRegarde 2011-12-20	352 × 352	327 × 327	681	397	332	30
LCP CaVousRegarde 2012-01-19	306 × 306	289 × 289	638	371	290	22
LCP CaVousRegarde 2012-01-25	356 × 356	350 × 350	747	428	352	25
LCP EntreLesLignes 2011-12-16	428 × 428	251 × 251	788	560	253	15
LCP EntreLesLignes 2012-01-27	417 × 417	264 × 264	840	599	267	14
LCP EntreLesLignes 2012-05-11	405 × 405	300 × 300	949	960	305	13
LCP LCPIInfo13h30 2012-01-24	159 × 159	158 × 158	337	212	163	13
LCP LCPIInfo13h30 2012-01-25	262 × 262	168 × 168	460	331	172	10
LCP LCPIInfo13h30 2012-01-27	225 × 225	170 × 170	456	306	177	11
LCP PileEtFace 2011-11-19	207 × 207	210 × 210	441	243	211	40
LCP PileEtFace 2011-12-01	293 × 293	254 × 254	587	389	255	51
LCP PileEtFace 2012-01-12	277 × 277	268 × 268	580	330	269	52
LCP PileEtFace 2012-01-19	281 × 281	244 × 244	518	302	246	43
LCP PileEtFace 2012-01-26	283 × 283	270 × 270	624	380	273	50
LCP TopQuestions 2012-01-25	258 × 258	130 × 130	480	364	134	11
LCP TopQuestions 2012-02-14	172 × 172	153 × 153	395	249	155	13
LCP TopQuestions 2012-02-22	214 × 214	155 × 155	439	308	158	15

### 3. RELATED WORK

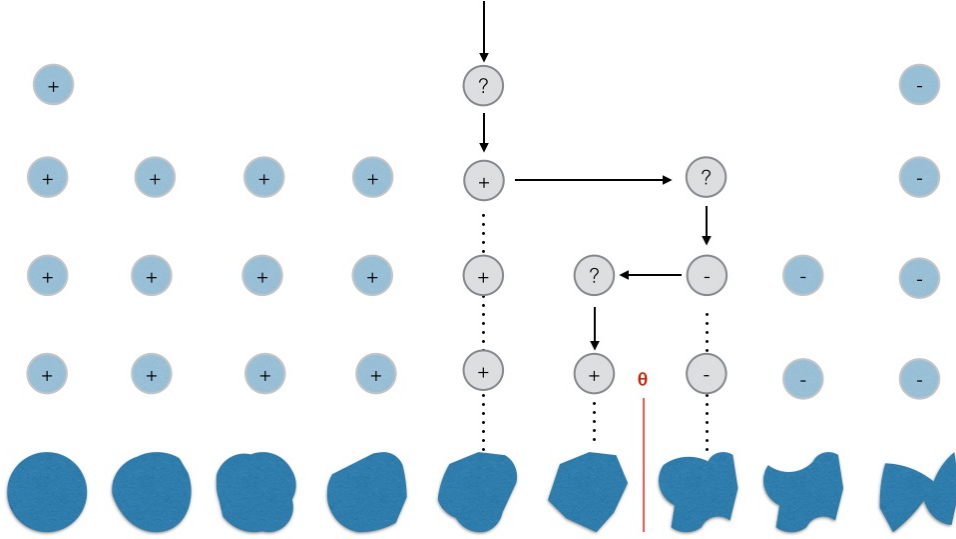
In this chapter, we review some fundamental concepts for this study. We give a brief overview of active learning, followed by a general discussion on active learning scenarios and basic active learning strategies.

#### 3.1 Active Learning

Active learning mainly discovers the most informative instances in terms of data distribution in data, while training a model in an iterative manner to deal with the cost of labeling. Determining the most informative instances and querying human annotators iteratively decreases the needed number of labeled data. Moreover, it helps the learner algorithm for learning the underneath distribution relatively faster than supervised methods. A nice example [14] about the human colonial, which arrives a solar planet, explains the semantic behind the active learning. In this scenario, the planet is habitable and includes a large amount of eatable vegetation. Important amount of the food comes from a plant of which fruits are smooth, round and irregular. The smooth and round fruits are delicious and good for humans. On the other hand, irregular fruits cause sickness. It is very indispensable to classify fruits as safe and noxious for the favor of the colony.

According to the PAC learning framework [15], if data distribution can be perfectly classified by some hypothesis function  $h$  in the hypothesis function set  $H$  then it is enough to test  $\mathcal{O}(1/\epsilon)$  randomly selected instances, where  $\epsilon$  is the maximum desired error rate. In other words, hundreds of fruits might be tested to be able to achieve 99% accuracy for 'fruit safety' classifier. Unfortunately, this experiment can cause lots of poisoned people. Conversely, instead of choosing the fruits randomly for specifying fruit irregularity threshold  $\theta$ , we can find and use the most informative fruits.

All things considered, we are able to formulate the classification problem of 'safe' and 'noxious' food by using simple binary classifier. Literally, each instance of the dataset



**Figure 3.1:** Detecting safe shaped food threshold  $\theta$  with binary search.

$X = x_1, x_2, \dots, x_n$  represents a harvested fruit. Each data instance  $x_m$  has a hidden label  $y \in \{1(\text{safe}), 0(\text{noxious})\}$ . Thus, the binary classification problem is defined as the set of  $x_m, y_m$  pairs with the function mapping  $h : X \rightarrow Y$ , parameterized by a threshold  $\theta$ , which represents the best threshold to be able to decide 'safe' and 'noxious' fruits. Supervised learning based methods need more human subjects to specify the threshold. Fortunately, binary search reduces the threshold search cost in this problem as seen in Figure 3.1.

The study [14] indicates that binary search speeds up the finding errorless classifier for a number of  $\mathcal{O}(\log_2 \frac{1}{\epsilon})$  training instances where  $\epsilon$  indicates the error rate. By this way, there is no need to use all of the training instances. For example, previous supervised methods use 100 human subjects to taste 100 fruits for obtaining 99% accuracy. On the other hand, the same accuracy can be achieved with 6 or 7 human subjects by using active learning based classifiers. In the end, using the most informative fruits enables a substantial reduction in the number of sick people.

### 3.1.1 Active Learning Scenarios

Active Learning scenarios defines how to ask a selected instance to a human expert. Mainly, there are 3 different Active Learning scenarios namely query synthesis, stream based selective sampling and pool based sampling which explained in the following sub sections in detail. In this study, we use stream based selective sampling.

### **3.1.1.1 Query Synthesis**

D. Angluin's 'Query and concept learning' [16] article describes the query synthesis. This scenario investigates learning with membership queries. In query synthesis learning scenario, learner may request label membership for any labeled data instance in the input space including queries that the learner starts from the beginning [17].

Query synthesis is good if learner has an exact definition of the input space. I realise in my literature search that query synthesis is often used. Especially recent works rarely use query synthesis. [14] gives a good example about query synthesis. [14] applies membership query learning with human oracles to train a neural network classifier of handwritten characters. The problem is many of query images generated by the learner can not be recognised by oracle. They encountered an unexpected problem: many of the query images generated by the learner contains no recognisable symbols; these are artificial hybrid characters with little or no natural semantic meaning for oracle. For example; hand writing tools can express 5, 8, or 9 instead of the real image of a number. It stands to reason that undetermined images could help the learner discriminate among the different characters, if people were able to discriminate among them as well. Similarly, text or speech recognition tasks are not suitable for query synthesis because of same sample generation problem.

### **3.1.1.2 Stream Based Selective Sampling**

Other query scenario is stream based selective sampling. Stream based selective algorithm is easy to understand. Learning starts with random labeled data, when the learning process, learner decides to ask one sample to oracle. The key assumption is that obtaining an unlabeled instance cost is none (or inexpensive), so it can be sampled from the actual distribution, and then the learner can decide whether or not to request its label. If the input distribution is uniform, selective sampling cannot offer any advantage over query synthesis. But if the input distribution is non-uniform, selective sampling gives better result than query synthesis [14].

### **3.1.1.3 Pool-Based Sampling**

Stream based selective sampling is needs more computation and time. Real world problems are not suitable stream based selective sampling. In real world problems, starting with small unlabeled data to training then learner ask oracle to large collections of unlabeled data. Learner creates a query pool when training [14].

The study [14] explains difference between stream based selective sampling and pool based sampling difference in his book. Main difference between stream based and pool based active learning is selection of query sample. Stream based query ask one sample at a time. But pool based method asks set of sample to expert after each learning iteration. Stream based query method looks more precise then pool-based because they ask one sample after each iteration this is why asked sample quality is guaranteed. But pool based sample can ask some noise or meaningless samples. On the other hand, stream based querying need more training time, because each annotation followed by training phase. This is why pool based querying is more popular if annotation cost is not too high unlike video annotation.

### **3.1.2 Sampling Strategies**

We explain how to ask scenarios in section two. Most important part of active learning is how to choose or decide which samples are best for the asking to oracle. I will explain three strategies, uncertainty sampling, query by disagreement, query by committee.

#### **3.1.2.1 Uncertainty Sampling**

Uncertainty sampling is proposed by [18]. This algorithm queried most confident samples for learner. Main idea is finding most informative instances from unlabeled input. So we can find most informative instances near the decision boundaries. If learner chooses instances near of decision boundary, probability near 0.5 for two class problem near decision boundary, learner can ask these confused points to oracle for labeling, so learner can train better in next iteration in theory.

The book [14] refers three measures of uncertainty in his active learning book.

- **Least Confident**

Least confident method is a basic uncertainty sampling strategy to query the instance whose predicted output is the least confident:

$$x^*_{CL} = \operatorname{argmin}_x P_\theta(\hat{y}|x)$$

$$= \operatorname{argmax}_x 1 - P_\theta(\hat{y}|x) \quad (1).$$

With the highest posterior probability under the model  $j$ . In other words, this strategy prefers the instance whose most likely labeling is actually the least likely among the unlabeled instances available for querying [14].

- **Margin**

A different active learning strategy is based on the output margin:

$$x^*_M = \operatorname{argmin}_x [P_\theta(\hat{y}_1|x) - P_\theta(\hat{y}_2|x)] \quad (2)$$

where  $\hat{y}_1$  and  $\hat{y}_2$  are the first and second most likely predictions under the model, respectively Margin sampling addresses a shortcoming of the least confident strategy by incorporating the second best labeling in its assessment [14].

- **Entropy**

Perhaps the most general (and most common) uncertainty sampling strategy uses entropy [19], usually denoted by  $H$ , as the utility measure:

$$x^*_H = \operatorname{argmax}_x H_\theta(Y|x)$$

$$x^*_H = \operatorname{argmax}_x - \sum_y P_\theta(y|x) \log P_\theta(y|x) \quad (3)$$

Where “ $y$ ” ranges over all possible labels of “ $x$ ”. Entropy is a measure of a variable’s average information content. As such, it is often thought of as an uncertainty or impurity measure in machine learning.

### 3.1.2.2 Query by Disagreement

Query by disagreement (QBD) proposed by [20]. Disagreement assumes the stream-based selective sampling scenarios so unlabeled data streams to oracle for labeling. Main idea is setting up two hypotheses; one hypothesis is most general hypothesis on the other hand other hypothesis sets up most specific hypothesis. Two classifiers work together but if one of classifier is disagree with other classifier instance is asked to oracle for labeling.

The study [14] indicates, QBD needs perfect approximate to boundaries for measure disagreement among all hypotheses in the version space, even if we imperfectly approximate it with two extreme hypothesis, most general and most specific QBD fall in trouble. Query by committee relaxes this assumption using multi classifiers.

### 3.1.2.3 Query by Committee

More theoretically-motivated query selection framework is the query by committee (QBC) algorithm [21]. Query by committee approach foundation is set of different classifiers. In this approach, set of different classifiers train for the same data set. In labelling section all classifiers propose own label to same sample. Same sample label agreed by all or most of classifiers but some sample cause conflict between classifiers. This disagreed samples are most informative data. Because those samples are near decision boundary in theory and we know that decision boundary is the most informative section of hypothesis space [14]. Bagging like method works well in query by committee [22].

Disagreement of classifiers need a measure method. Two dominant methods used for the measuring disagreement in classification task. First one is based on entropy, which called vote entropy. This measurement is most used measurement in literature. Formulation is given [14]:

$$x_{VE}^* = \underset{x}{\operatorname{argmax}} \sum_y \frac{\operatorname{vote}_c(y,x)}{|C|} \log \frac{\operatorname{vote}_c(y,x)}{|C|}$$

$$\operatorname{vote}_c(y,x) = \sum_{\theta \in C} 1_{h_\theta(x)=y} \quad (4)$$

In this formulation,  $y$  is all possible labels,  $|C|$  size of the committee.



Another disagreement measure is Kullback-Leiber (KL) [23] . This method tests measurement of the difference between two probability distributions. In this case , KL quantify disagreement as the average divergence of each committee member  $\theta$  prediction from that of the consensus C. Also formulation is given as the following;

$$x_{KL}^* = \underset{x}{\operatorname{argmax}} \frac{1}{|C|} \sum_{\theta \in C} KL(P_{\theta}(Y|x) || P_C(Y|x))$$

$$KL(P_{\theta}(Y|x) || P_C(Y|x)) = \sum_y P_{\theta}(y|x) \log \frac{P_{\theta}(y|x)}{P_C(y|x)} \quad (5)$$

[14] mentions about KL that some cases probability distributions is not uniform even if not consensus output may be uniform. This is difference between vote entropy and KL divergence.

### 3.2 Cluster Matching Metrics

[24] proposes theoretic "Variation Of Information"(VI) criterion for comparison of two clustering on the same data set. Proposed VI criterion measures and compares the qualities of two different C and C' clusters to determine the better clustering method in terms of information gain and loss. In addition, the study gives a nice literature review about the comparison methods of previous studies and groups them. The first group of comparison methods uses Counting Pairs of compared cluster methods. However, Counting Pair methods do not give better results than other comparison metrics because of the asymmetry. Pair of points from X described below for all cases:

- $S_{11}$  demonstrates the number of the data point pairs that are in the same cluster under both C and C'
- $S_{00}$  demonstrates the number of the data point pairs that are in the different clusters under both C and C'
- $S_{10}$  demonstrates the number of the data point pairs in the same cluster under C, but in different clusters under C'
- $S_{01}$  demonstrates the number of the data point pairs in the same cluster under C', but in different clusters under C

The sum of all the cases always must be;

$$S_{11} + S_{00} + S_{10} + S_{01} = n(n-1)/2 \quad (3.1)$$

By considering the Eq.3.1, [25] proposes two asymmetric comparison criteria Eq.3.2 and Eq.3.3 which gives the probability of point pairs in the same cluster under both  $C$  and  $C'$ .

$$W_I(C, C') = \frac{S_{11}}{\sum_i n_i(n_i - 1)/2} \quad (3.2)$$

$$W_{II}(C, C') = \frac{S_{11}}{\sum_j n'_j(n'_j - 1)/2} \quad (3.3)$$

[26] proposes a comparison metric called Rand Index which compares the result of a classification scheme with a correct classification (See Eq.3.4 ). Rand Index depends the number of clusters and the number of data points. In the equation, the function  $R$  gets values between [0,1] where  $R = 0$  indicates that the two clustering algorithms are identical.

$$R(C, C') = \frac{2(n_{11} + n_{00})}{n(n-1)} \quad (3.4)$$

[24] mentions another criteria called Jaccard Index as in Eq.3.5. Jaccard Index gets values between [0,1] where  $J = 0$  indicates that the two clustering algorithms are identical.

$$J(C, C') = \frac{S_{11}}{S_{11} + S_{01} + S_{10}} \quad (3.5)$$

[27] improves the Jaccard Index as in Eq.3.6

$$M(C, C') = \sum_i n_i^2 + \sum_j n'_j{}^2 - 2 \sum_i \sum_j m_{ij}^2 \quad (3.6)$$

[24] mentions Set Matching Criterion family as the second group of comparison methods which compare clusters using cluster sets not data points. The first Set Matching Criterion called Larsen and Aones criterion which is asymmetric

(See Eq.3.7). The criterion function  $L$  gets values between [0,1] where  $L = 1$  indicates that the two clustering algorithms are identical.

$$L(C, C') = \frac{1}{k} \sum_i \max_j \frac{2m_{ij}}{n_i + n'_j} \quad (3.7)$$

[28] proposes criterion in the Eq.3.8 which finds the best matches of each cluster  $C$  in  $C'$  and is symmetric. It scans element of  $m_{ij}$  of the contingency table in decreasing order. Largest elements of  $m_{ij}$  called as nab. All row and column matches are listed until to reach  $\min(k, \ell)$ .  $match(i)$  function is index of  $C'_i$  in  $C'$  that matches cluster  $C_i$ . The criterion function  $H$  gets values between [0,1] where  $H = 1$  indicates that the two clustering algorithms are identical.

$$H(C, C') = \frac{1}{k} \sum_{j=match(i)} m_{ij} \quad (3.8)$$

Another symmetric criterion introduced by [29] is given in Eq.3.9.  $D = 0$  indicates that the two clustering algorithms are identical, otherwise  $D$  value is smaller than  $2n$ .

$$D(C, C') = 2n - \sum_i^k \max_j m_{ij} - \sum_j^\ell \max_i m_{ij} \quad (3.9)$$

All the three matching symmetric method functions  $L$  in Eq.3.7,  $H$  in Eq.3.8 and  $D$  in Eq.3.9 finds the best matching sets first for each cluster and compares the clusters using best matching values. The study [24] mentions the problem of these three criterion do not compute "unmatched" part of each cluster. So these metrics do not know what happens clusters unmatched points. It introduces VI criterion for solving asymmetric and unmatching problem.

Variation of information is aimed to compare clusters  $C$  and  $C'$  using disagreement. It uses entropy based method for measuring disagreement. First step of finding entropy of cluster is calculate probability of cluster sets:

$$P(i) = \frac{n_i}{n} \quad (3.10)$$

Thus, we can calculate the entropy of each cluster as;

$$H(C) = \sum_1^k P(i) \log P(i) \quad (3.11)$$

The formula in Eq.3.11 shows the entropy measured in bits and entropy does not depend of number of points. this calculation is basement to calculation of mutual information between two clustering.

Calculating of mutual information need  $P(i, j)$  this is joint distribution of the random variables.

$$P(i, j) = \frac{|C_i \cap C'_j|}{n} \quad (3.12)$$

Mutual information of associated random variables described in below;

$$P(i, j) = \sum_{i=1}^k \sum_{j=1}^{\ell} P(i, j) \log \frac{P(i, j)}{P(i)P'(j)} \quad (3.13)$$

[24] introduces variation of information using mutual information and entropy associated with cluster C.

$$VI(C, C') = [H(C) - I(C, C')] + [H(C') - I(C, C')] \quad (3.14)$$

We can see that  $I(C, C')$  is intersection of  $H(C)$  and  $H(C')$  this is why [24] proposed criterion find disagreement of two cluster subtracting mutual information from each entropy associated cluster. It examined VI under 12 property like non-negativity, symmetry, triangle inequality, n-invariance, upper bound of VI, splitting a cluster, collinearity of join, linearity of composition.

Another study about evaluation metrics to compare clustering algorithms done by [30]. Amigo et al. defines four new formal measures for clustering comparison, writers also examines some comparison metrics using newly defined constraints. Amigo et al. also explain why those metrics can not satisfy all constraints in first part of study. Second part of study Authors analyze BCube algorithm because BCube satisfies all of four constraints also BCube can compare clustering that has overlapping problem.

[30] tested some comparison metrics - explained before- using 12 mathematical constraints. It defines higher level formal constraints than [24] constrains because authors expect some rules for constraints:

A Constraint should be intuitive and tests clearly limitation of each metric. So system developer can be sure which constraints are important for the specified task. Constraint should be prove formally. Metric families must be discriminated by constraints. This property is useful for system developer because system developer can choose right metric if previous tested family does not work properly for the problem.

Using this point of view, authors define those formal constraints on evaluation metrics for clusterings.

Constraint one is cluster homogeneity: This constraint is purpose of clustering algorithms, clusters must be homogeneous. In a word, all clusters have same item. Second constraint is cluster completeness: First constraint measure homogeneity but same class items can be clustered in two or more cluster. In many case, big one cluster is better than two or more clusters for same class items. Third constraint is rag bag . Rag bag is clustering noise points to same cluster that means rag bag is set of 'others', 'unclassified', 'miscellaneous' categories. Constraint four is cluster size versus quantity. Big cluster and small error should be preferable to large number of small error in small clusters. For example if two clustering algorithms has same error number for same data, clustering algorithm that has bigger cluster is better than other algorithm.

Authors examine families of metrics under those constraints, results:

Set matching metric families can not satisfy constraints 2 and 3. Constraint 1 is satisfied by purity measure. Counting pair based metrics family-like Rand, Jaccard Coefficient, Folkes and Mallows can satisfy constraints 1 and 2 but they do not satisfy 3 and 4. Entropy based metric family like Meila's Variation of information fail to satisfy constraints 2 and 4. Evaluation metrics family based on edit distance can not satisfy constraints 1 and 3. [30] introduce BCubed metrics - metric that is mixed family of metrics- can satisfy all of four constraints.



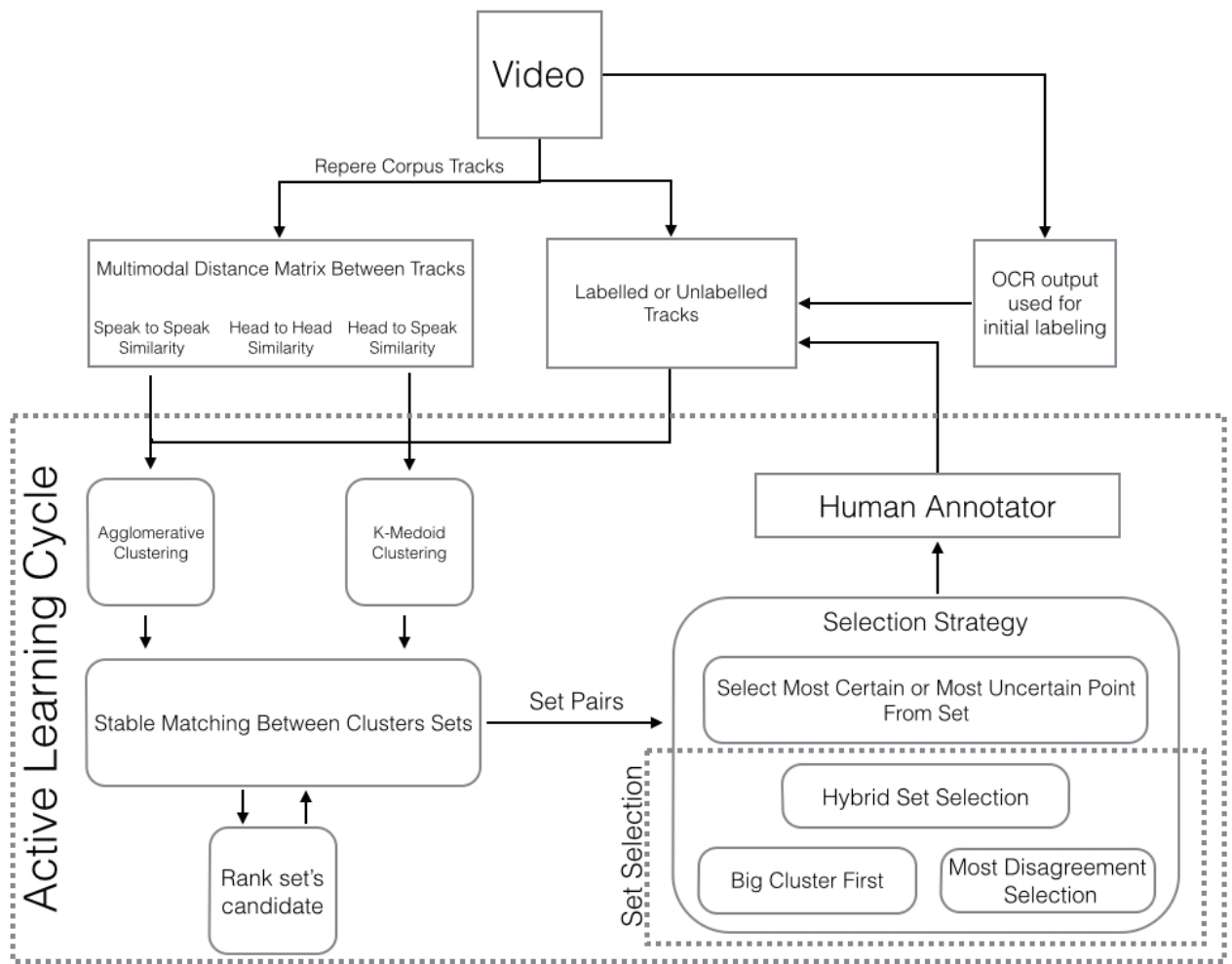
## 4. METHODOLOGY

Supervised query by disagreement method (QBD) [14] is one of the active learning methods in the literature. QBD uses two different classifiers and asks the label of a disagreed point to an expert. We introduce unsupervised query by disagreement method to accelerate the learning phase of the automatic video annotation by using unsupervised learners. When it is used with supervised learners, Query by Disagreement (QBD) lets the learners label each instance and compares the learners' outputs for the same instance to detect disagreement instances. However, since unsupervised methods do not generate labels for each data instance as output of a learner, using QBD on unsupervised learners is a new and a challenging problem. We propose a novel approach for QBD on unsupervised learners and apply the solution on multimodal video data.

In our study, we use the following steps during an Active Learning Cycle: clustering, cluster matching, disagreement measurement between clusters, selection of the most disagreed clusters, selection of an instance to be queried from the selected cluster. The architecture of our system is depicted in Figure 4.1

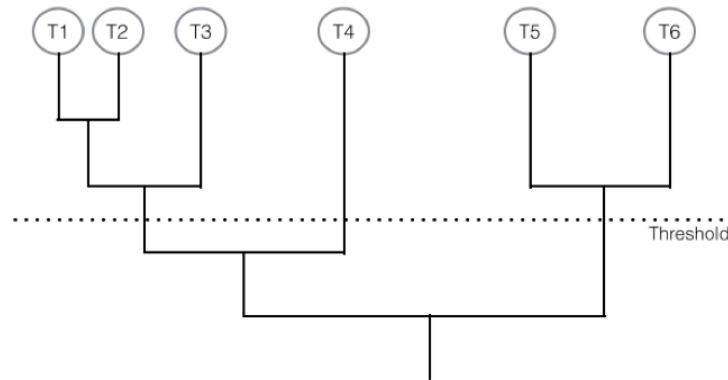
### 4.1 Clustering Algorithms

Active learning cycle begins with the clustering of the data. We use two clustering algorithms: Agglomerative Clustering and K-Medoid Clustering. Agglomerative clustering is a hierarchical clustering algorithm and works with bottom-up approach. Therefore, it does not need to take a fixed number of clusters as input. For a given distance threshold, Agglomerative Clustering estimates the number of clusters, which we use as the value of the K-Medoid Clustering's number of clusters parameter  $k$ .



**Figure 4.1:** The architecture of the system.





**Figure 4.2:** The illustration of Agglomerative Clustering.

#### 4.1.1 Agglomerative Clustering

Agglomerative Clustering method clusters instances up until reaching a given distance threshold and returns the number of clusters. In our dataset, each video data has variable size of classes. However, the other clustering algorithm that we use, K-Medoid, needs number of clusters as input. Thus, we let the Agglomerative Clustering finds the number of the clusters and use it as input in K-Medoid Clustering. The generic Agglomerative Clustering algorithm is given in the Figure 4.3

```

Initialize
Assign each data point as a cluster
while until reaching a given distance threshold do
    Find the most similar pair of clusters
    Merge them into a single cluster
    Compute new similarities of the new clusters
end while

```

**Figure 4.3:** Pseudo-code of Agglomerative Clustering.

#### 4.1.2 K-Medoid Clustering

K-Medoid Clustering algorithm needs an initial  $k$  parameter which corresponds to the number of clusters. However, our proposed unsupervised query by disagreement needs at least one clustering method that does not need to know the number of clusters initially. Therefore, we use agglomerative clustering to gather the number clusters and use it as K-Medoid's parameter  $k$ .

```

Initialize
k points are randomly selected in n instances as medoids
while Is medoids shifted do
  Each data point is associated to the closest medoids
  for each medoid m do
    for each datapoint 'o' is not a medoid do
      swap m and o and calculate the total costs of the configuration
    end for
  end for
end while

```

**Figure 4.4:** Pseudo-code of K-Medoid algorithm.

In order to find the clustering disagreements between the algorithms, we first solve its complementarity problem by finding the agreed instances on correspondent clusters the of algorithms. We apply a cluster matching algorithm to measure the similarities between the clusterings produced by the Agglomerative and K-Medoid algorithms. We find the intersection of clusters by calculating the cluster similarity metric given in Equation 3.8. For our dataset, the cluster pairs' similarity matrix using best match measurement is given in Table 4.1.

**Table 4.1:** A sample cluster pair similarity matrix for agglomerative and K-Medoid algorithms.

		<b>K-Medoid</b>		
		<b>Cluster 1</b>	<b>Cluster 2</b>	<b>Cluster 3</b>
<b>Agglomerative</b>	<b>Cluster 1</b>	0.437	0.094	0.321
	<b>Cluster 2</b>	0.318	0.432	0.176
	<b>Cluster 3</b>	0.134	0.126	0.173
	<b>Cluster 4</b>	0.320	0.284	0.421

After finding the cluster pairs' matching similarities, we need to couple the pairs using cluster pair similarity matrix that we constructed using 'Best Match'. We calculate the similarity matrix between the sets of two clustering algorithms but we do not assign Agglomerative Clustering's sets to K-Medoid Clustering's sets one by one yet. The 'assigning a set to another set' problem perfectly fits into 'Stable Marriage Problem'. In real life, finding the most coherent man and woman pair is the key of a stable

marriage. However,  $n$  men and  $n$  women stable marriage is an optimization problem because all candidates are assigned to the most coherent candidates. Each man and woman ranks all of the opposite members from 1 to  $n$  and a stable marriage occurs when a man and a woman could not find any other better partners in opposite sex within  $n - 1$  candidate in terms of rankings.

Literally, we need to find  $j = match(i)$  which denotes the cluster  $C'_j$  to which the  $C_i$  should be matched. Then, in order to assign clusters between Agglomerative and K-Medoid clusterings, we adopt and apply the Gale-Shapley algorithm ([31], [32]) which solves the 'Stable Marriage Problem'. Gale-Shapley algorithm guarantees that the solution obtained is perfect (i.e. everyone gets married) and stable.

Let the clustering produced by the Agglomerative Clustering be  $C = C_1, C_2, \dots, C_k$  and the Clustering produced by the K-Medoid Clustering be  $C' = C'_1, C'_2, \dots, C'_k$ . Gale-Shapley algorithm requires that each cluster  $C_i$  ranks the clusters  $C'_j$  and vice versa. We use the cluster confusion matrix entries (Equation 2.2) in order to produce these rankings and then apply the Gale-Shapley algorithm to produce a cluster matching. The preferences of each Agglomerative cluster over K-Medoid clusters are denoted as  $P$ . The preference of the cluster  $C_1$  may be shown as  $P(C_1) = C'_2, C'_4, C'_1, \dots$  where each cluster in the list  $P(C_1)$  comes from the opposite cluster set  $C'$  and in descending order in terms of the candidates' scores. In the same manner, the preferences of each K-Medoid cluster over agglomerative cluster set  $C$  are listed. The preference list of the whole agglomerative cluster set is denoted as  $P_a = P(C_1), P(C_2), \dots, P(C_i)$ . Likewise, the preference list of the whole K-Medoid cluster set is denoted as  $P_k$ . The output of the Gale-Shapley is a matching  $M = [m_1, m_2, \dots, m_k]$  where,  $m_i \in 1, \dots, k$ ,  $m_i = j'$  if  $C_i$  is matched with  $C'_j$  in the stable matching. Stable matching algorithm for two clusters  $C$  and  $C'$  is given in Figure 4.5.

## 4.2 Query Instance Selection Strategies

In order to be able to select which instance to query, we first select the most informative cluster of the cluster pair, and then we select the most informative instance in the selected cluster.

```

for each cluster  $C_i \in C$  do
   $P(C) \leftarrow \text{BestMatch}(C_i, C')$ 
  Add  $P(C)$  into  $P_a$ 
end for
for each cluster  $C'_i \in C'$  do
   $P(C') \leftarrow \text{BestMatch}(C'_i, C)$ 
  Add  $P(C')$  into  $P_k$ 
end for
while unstable do
  unstable  $\leftarrow$  false
  for each cluster  $C_i \in C$  do
    if  $C_i$  has a pair cluster in  $M$  or  $P(C_i) = \emptyset$  then
      continue
    end if
    unstable  $\leftarrow$  true
     $C'_i \leftarrow$  the first element of  $P(C_i)$ 
    if  $C_i$  has not a pair cluster in  $M$  then
       $m_x \leftarrow \langle C_i, C'_j \rangle$ 
      Add  $m_x$  into  $M$ 
       $m_y \leftarrow \langle C'_j, C_i \rangle$ 
      Add  $m_y$  into  $M$ 
    else
       $m \leftarrow$  current match for  $C'_i$ 
      if Rank of  $\langle C'_j, C_i \rangle >$  Rank of  $m$  then
         $m_x \leftarrow \langle C_i, C'_j \rangle$ 
        Add  $m_x$  into  $M$ 
         $m_y \leftarrow \langle C'_j, C_i \rangle$ 
        Add  $m_y$  into  $M$ 
        Delete current match  $m$  from  $M$ 
      end if
    end if
  end for
end while

```

**Figure 4.5:** The stable matching algorithm.

## 4.2.1 Cluster Selection

We score clusters using different methods and select an instance from the cluster with the highest score.

### 4.2.1.1 Big Cluster First

The study [2] proposes 'Big Cluster First' (BCS) selection strategy which calculates a score using the size of a set and the number of annotated instances in that set. The method selects an instance from a minimum scored set by asking human expert. The BCS strategy score for a cluster  $C_i$  is calculated as shown in Equation 4.1.

$$BCS(C_i) = \frac{NumberOfAnnotation(C_i)}{Size(C_i)} \quad (4.1)$$

### 4.2.1.2 The Most Disagree Selection

In theory, the most disagreed cluster pair gives us the most uncertain points because two stable matched clusters have lots of disagreed instances. MDS method chooses the cluster  $C_i$  which has the highest DS score. For measuring the disagreement between a pair of matched clusters  $C_i, C'_j$  where  $j = m_i$ , we calculate the disagreement score (DS), disagreement score list (DSL) and most disagree pair (MDS) as the following;

$$m_x = \langle C_i, C'_j \rangle, \forall m \in M \quad (4.2)$$

$$DS(\langle C_i, C'_j \rangle) = \frac{C_i \cap C'_j}{C_i \cup C'_j} \quad (4.3)$$

$$DSL \{ DS(\langle C_i, C'_j \rangle) | \forall m \in M : m_x = \langle C_i, C'_j \rangle \} \quad (4.4)$$

$$MDS = (\langle C_i, C'_j \rangle), \text{ where } DS_s(\langle C_i, C'_j \rangle) = \min(DSL) \quad (4.5)$$

where

$$DS_s(\langle C_i, C'_j \rangle) = \min(DSL) \quad (4.6)$$

### 4.2.1.3 Hybrid Cluster Selection

BCS method performs better during the initial stages of active learning where the number of labeled instances are very few and labeling instances on big clusters help label a lot of instances. On the other hand, MDS performs better when clusters contain more known instances, because cluster disagreement has a correlation with the label assignment. Therefore, we introduce two hybrid cluster selection methods: Soft Hybrid Selection (SH) and Hard Hybrid Selection (HH).

In order to combine or compare the different BCS and MDS clustering scores, we apply the z-score normalization on them.

We know *BCS* performs better than random selection for early iterations in active learning cycle. On the other hand, *MDS* is not good at early stages but *MDS* performs good after few iterations later. Because of this manner, we develop the **Soft Hybrid Selection** method. Contracting SH method needs *BCS* like method for cluster pairs. This is why we propose big cluster pair first selection method for use in early iterations. The Big Cluster Pair First method is the derivative of *BCS* for paired clusters. The *BCPS* method calculates a score for each pair and selects lowest scored pair.

A cluster pair score is depicted as  $S_p = A_p/N_p$  where  $A_p$  is the total number of annotation in a cluster pair and  $N_p$  is the total number of instances in a paired cluster. Soft hybrid method combines *BCPS* and *MDS* scores using weight. However, *MDS* and *BCPS* scores' mean and standard deviations are not comparable since they use different scales. Therefore, z-score normalization is applied over both *BCPS* and *MDS* scores for solving the scale problem. After that we calculate SH score (SHS) for each matched pair score using the given formula 4.7. Normalized *BCEP* and *MDS* scores are denoted as  $S_n$  and  $DS_n$ , respectively. The soft hybrid score *SHS* is calculated as;

$$SHS = ((1 - \alpha) \times S_n) + (\alpha \times DS_n), 0 < \alpha < 1 \quad (4.7)$$

This formula clearly shows that if the value of  $\alpha$  is close to 1, selection shifts to *MDS* while lower  $\alpha$  value shifts SH selection to *BCPS*.  $\alpha$  value must be close to zero at early stage of active learning for selecting cluster pair from *BCPS*. Calculation of  $\alpha$  value is dependent to number of the iterations.

The **Hard Hybrid Selection** method uses *BCS* method for the given number of iteration. Important part is *BCS* uses only agglomerative clustering clusters. When active learning strategy cycle reaches to switching point, algorithm starts to use *MDS* for selection. In other words, hard hybrid uses *BCS* for early stages and switches to *MDS* later.

#### **4.2.2 Instance Selecting Strategy**

Selecting the most informative clusters with the defined cluster selection methods is the first step of the selection strategy. The next step is determining the most informative instances in the most informative clusters. The instances from the center of a cluster are more 'certain' than the instances close to the cluster boundary in terms of class knowledge. However, since the entropies of uncertain instances are relatively higher than the entropies of certain instances [14], the most informative instances are actually on the region around the boundary of a cluster. Therefore, for each instance in a cluster, we add the distance to other instances in the same cluster and choose the instance with the highest sum, which gives us the most 'uncertain' instance, the instance which is the farthest away from all the other instances. Instance selection strategy can select the most 'uncertain' or the most 'certain' instance from a cluster or a cluster pair.





## 5. EXPERIMENTS AND RESULTS

We evaluate the performance of selection strategies over active learning steps on REPERE Corpus data.

### 5.1 Evaluation Metrics

The confusion matrix also known as contingency matrix is basic metric of measurement of classifiers performance. The confusion matrix is an error for visualisation of classifier performance. The confusion matrix has class size of rows and columns. Each row contains total number of actual class value. Each column gives number of classifiers prediction. The confusion matrix reports four different data, *truepositives*, *falsepositives*, *truenegatives*, *truenegatives*. True positive indicates both actual class and the predicted class are positive. On the contrary, true negative indicates both actual class and predicted class are negative. On the other hand, false positive occurs when the actual class is negative but predicted class is positive. Likewise, false negative indicates the actual class is positive but predicted class is negative.

In this study, we use F-measure, which is an indicator of test set accuracy and harmonic mean of precision  $p$  and recall  $r$  values, to evaluate the algorithms. Precision shows how much relevant instances are retrieved among all retrieved instances. Regardless, Recall shows how much retrieved instances are relevant among all relevant instances. F-measure is calculated as in the Equation 5.1

$$F - Measure = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (5.1)$$

where

$$Precision = \frac{TP}{TP + FP} \quad (5.2)$$

and

$$Precision = \frac{TP}{TP + FN} \quad (5.3)$$

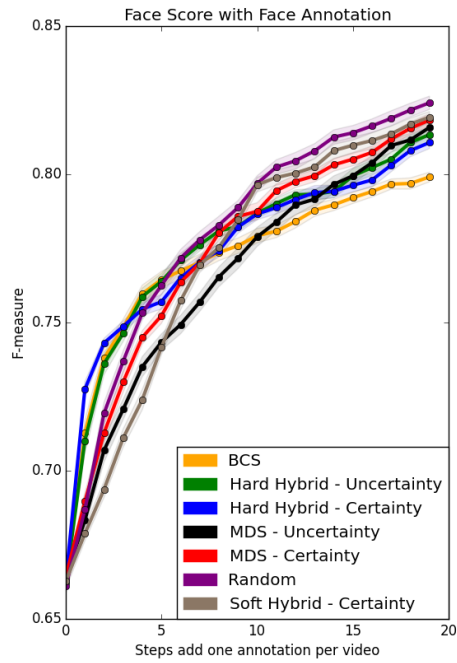
## 5.2 Results

We have evaluated the performance of selection strategies over active learning cycles using multimodality on 28 videos from 7 different TV programs. Each video has three different similarity matrices namely face-to-face, speech-to-speech and face-to-speech. We have run four experiments with regard to matrices in this order: (1) face score for face track annotation (FF), (2) face score for speaker track annotation (FS), (3) speaker score for speaker track annotation (SF), and (4) speaker score for speaker annotation (SS). An active learning cycle, which is depicted as *step*, asks one annotation for each video. As the performance measure, we use the F-measure, instead of accuracy, since the number of instances of each class (person) in the datasets is very different.

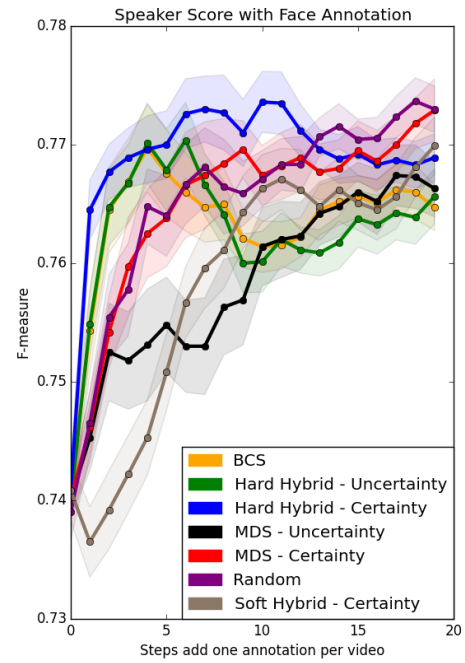
We have designed two different evaluation techniques to measure the performance of the algorithms. The first evaluation technique calculates overall performance for all videos by admitting 28 videos as one big video combining speaker and head tracks from all videos. We have then calculated FF, FS, SF, and SS scores (Figure 5.1, 5.2, 5.3 and 5.4).

The multimodal classifying problems FS and SF as seen in Figure 5.2 and 5.3 have more deviations than the others because the correlation between face tracks to speaker tracks is weaker than head to head tracks or speaker to speaker tracks. Therefore, in FS, MDS-Certainty performs better than all others, and more robust than BCS.

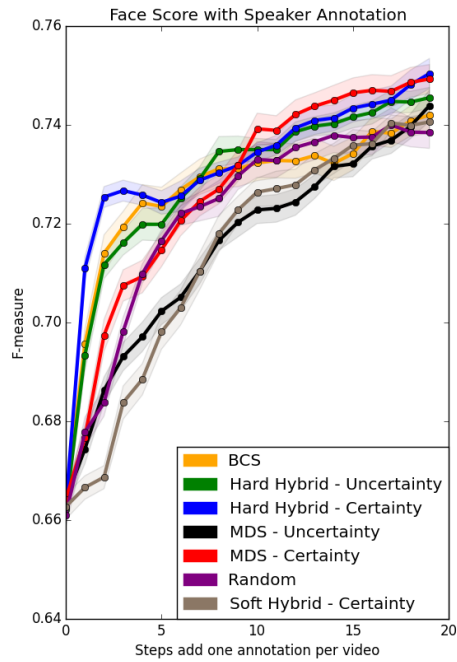
In the FS experiment, BCS selects good annotations at the earlier steps, and achieves high F-measure value faster than MDS-Certainty since there is no sufficient data early on to be disagreed on the MDS-Certainty method. However, at the later steps, F-measure value decreases, but MDS-Certainty keeps increasing. To take advantage of



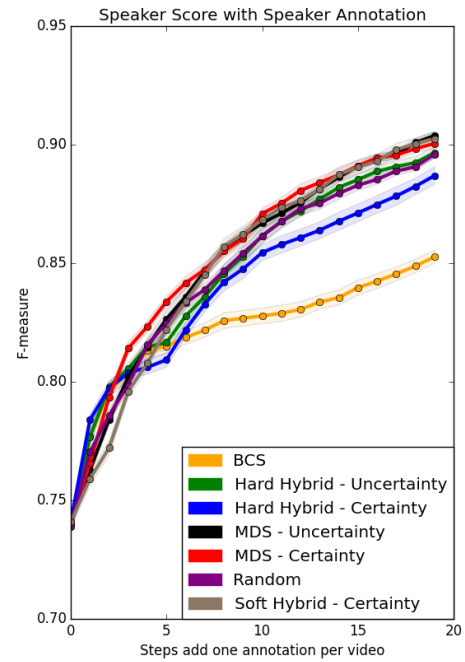
**Figure 5.1:** Face score with face annotation.



**Figure 5.2:** Speaker score with face annotation.



**Figure 5.3:** Face score with with speaker annotation.



**Figure 5.4:** Speaker score with with speaker annotation.

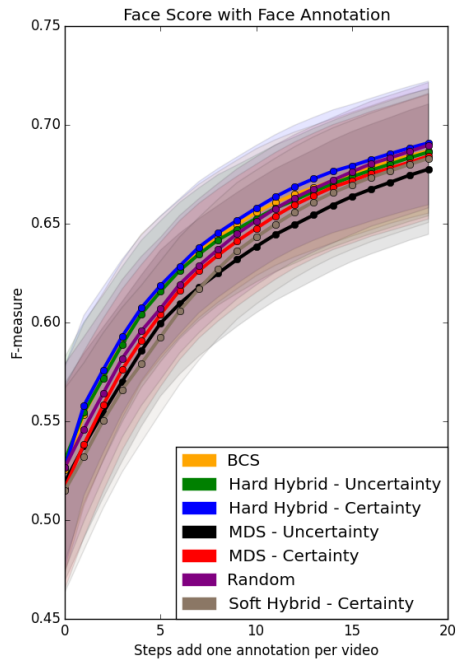
the strongest sides of both methods, Hard Hybrid Certainty (HHC) uses BCS at the first five steps. After the fifth step, it uses MDS and its F-measure score keeps increasing. As a result, HHC achieves better scores than random.

In the SF experiment, BCS achieves high F-measure scores rapidly and higher than random method like in FS. Nevertheless, it decreases at the later steps, and finally performs worse than random method. On the other hand, MDS-Certainty F-measure increases stably like that of FS, but it performs similar to random. However, the most interesting result comes from HHC that combines BCS and MDS-Certainty. Like BCS, HHC gives higher F-measures rapidly at the earlier steps; however, it continues to increase at the later stages like MDS-Certainty. Furthermore, it performs significantly better than random until the 10th step. Fortunately, for active learning, reaching to a high F-measure score rapidly is more important.

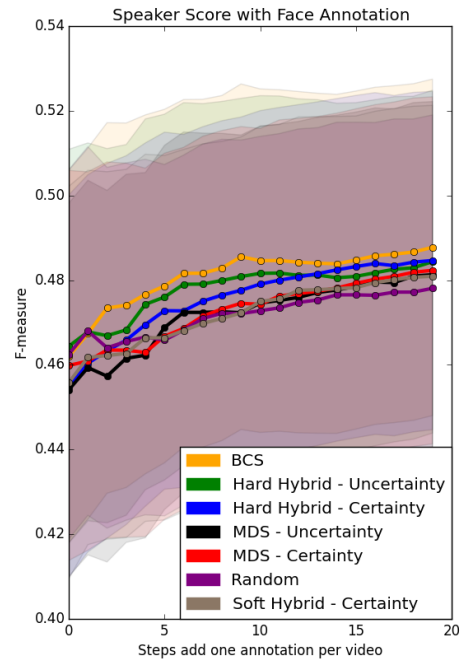
The second evaluation calculates F-measure scores for each video individually and returns one F-measure value that is the average of all F-measures. By doing so, we observe the performance of the algorithms video by video that have different size of similarity matrices. The video based F-measure results are given in ??for the four types of experiments.

We have noticed that some MDS based methods perform well when we use speaker annotation; especially HHC performs better for speaker identification with speaker annotation. We have calculated average values of video based performances. The results are given in Figure 5.5, 5.6, 5.7 and 5.8.

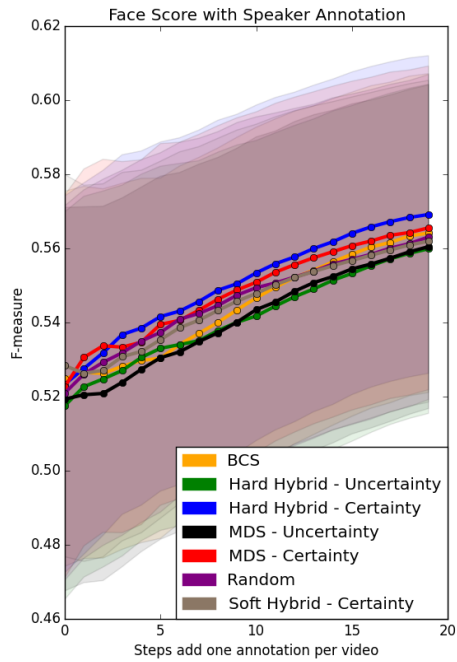
FF result in Figure 5.5 shows us HHC is barely better than other methods, and worse than BCS for SF experiment in Figure 5.6. On the other hand, HHC scores are better than speaker annotated experiments with FS and SS shown in Figure 5.7 and Figure 5.8. Another finding is that the deviations are high between video performances. For example, F-measure scores for some videos are between 0.7-0.9 while they are zero for some other videos. Such deviations indicate that we lose some information about the performance analysis of the methods using average video results. We have also ranked between active learning methods performance step by step for each videos, and created heat maps for all videos and all experiment types (Appendix A.1 ). We have calculated the average ranks using video based F-measure score averages, and average rankings of the methods for discussion since we apply active learning methods to different video types.



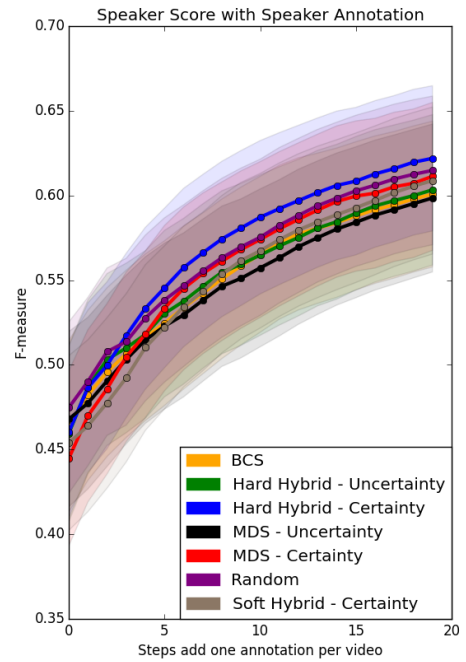
**Figure 5.5:** Face score with face annotation



**Figure 5.6:** Speaker score with face annotation.

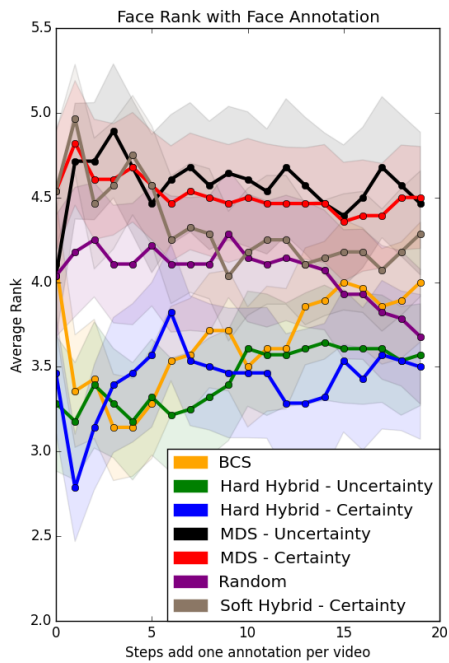


**Figure 5.7:** Face score with with speaker annotation.

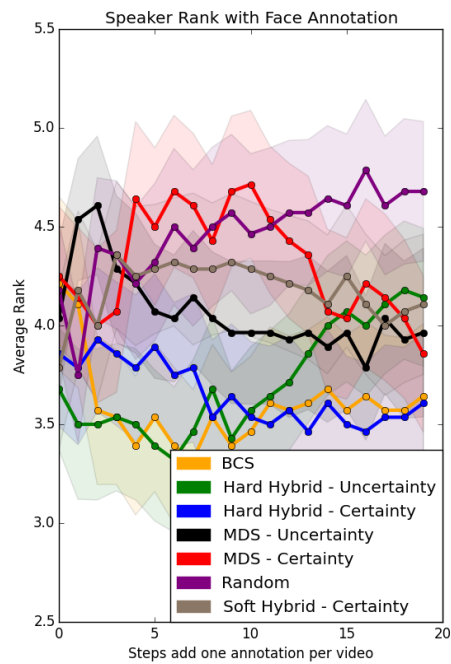


**Figure 5.8:** Speaker score with with speaker annotation.

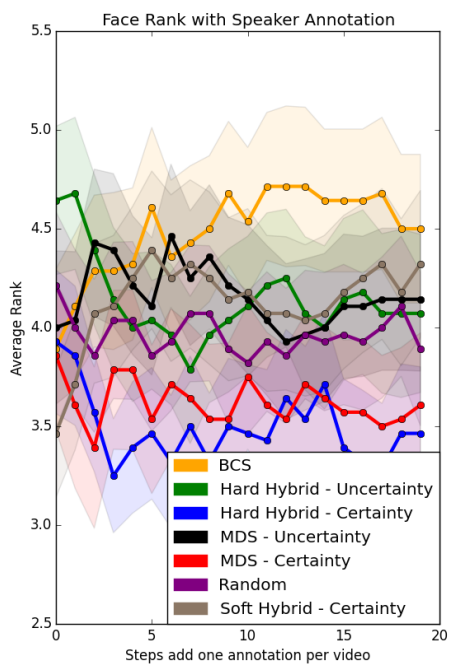
Our first finding is that HHC and HHU have the best performance rank average with BCS, but random selection beats BCS at late steps in FF experiment. On the other hand, both MDS-Certainty and MDS-Uncertainty performances are worse than



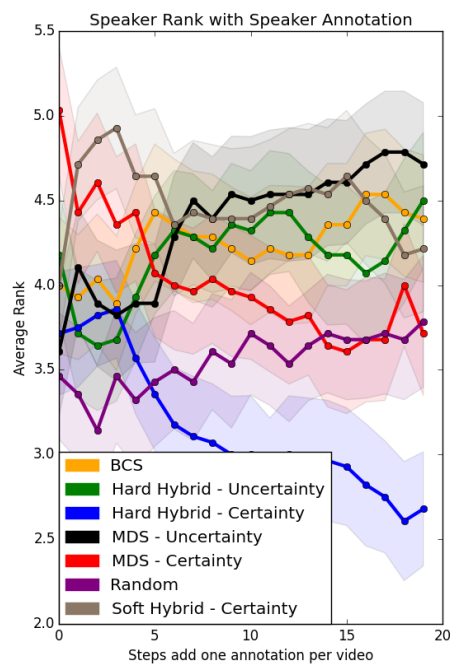
**Figure 5.9:** Face score with face annotation (lower is better).



**Figure 5.10:** Speaker score with face annotation (lower is better).



**Figure 5.11:** Face score with with speaker annotation (lower is better).



**Figure 5.12:** Speaker score with with speaker annotation (lower is better).

others. By examining the similarity matrices, we have found out that face-to-face average similarities are approximately ten times lower than speaker-to-speaker average similarities. Both MDS methods do not perform well because of lack of information about the data distribution for disagreement. Hard Hybrid (HH) method performs better due to BCS and MDS since at first 5 steps HH uses BCS, which collects sufficient information on early steps, and MDS which works better after the fifth step.

BCS has better average F-measure score, and ties with HHC in average ranking but all active learning methods beat random selection in SF experiment (Figure 5.10 and 5.10). On the contrary, HHC and MDS-Certainty outperform BCS and random in terms of both average F-measure score and average rankings (Figure 5.11 and 5.11). When we use cross modal annotation, total performance increments are less than single modality.

We know that speaker-to-speaker similarities higher than face-to-face similarities. HHC gives best performance in SS experiment (Figure 5.12) because of higher speaker similarities. We conclude that HHC is better than other sampling methods when similarities between annotated instances are high. High track similarity provides better clustering, and better clustering provides sharp disagreement or agreement. This is why MDS based hybrid method works well.





## 6. CONCLUSIONS AND RECOMMENDATIONS

Active learning mainly aims to reduce annotation/labeling cost. Active learning is a research area in machine learning and active learning is semi-supervised method. Main idea of the active learning is finding most informative instances for annotating or labeling train data. The most informative instance means that instance gives more information about data distribution so learner can learn data distribution quickly than supervised learning in theory. Annotation cost is important for some cases like video annotation because annotator spends long time and attention to labeling the video so this process is very costly and researchers research active learning algorithms to solving annotation cost problem. In addition, We have big data today especially Internet content like social networks, blogs, video sharing sites etc. annotation cost is high for many cases for big data because this is big, noisy, sparse this is why big data needs more annotation and active learning can reduce cost of big data annotation.

We proposed the MDS active learning method and its hybrid variations and we applied them on multimodal video annotation data. According to our experiments, for different types of annotation tasks, different active learning strategies could be more suitable. Hybrid strategies could be more successful than using a single strategy alone. Decision of the cluster selection and instance selection method adaptively during each active learning step, using a synthetic dataset to investigate the merits of these strategies, examination of each method for each video, rather than the whole REPERE corpus are the future research directions we aim to follow.



## REFERENCES

- [1] **Giraudel, A., Carré, M., Mapelli, V., Kahn, J., Galibert, O. and Quintard, L.**, 2012. The REPERE Corpus: a multimodal corpus for person recognition., LREC, pp.1102–1107.
- [2] **Budnik, M., Poignant, J., Besacier, L., Quénot, G. et al.**, 2014. Automatic propagation of manual annotations for multimodal person identification, Proceeding of the 12th International Workshop on Content-Based Multimedia Indexing.
- [3] YouTube Statistics 2014, <https://www.youtube.com/yt/press/statistics.html>.
- [4] **Ayache, S. and Quénot, G.**, 2007. Evaluation of active learning strategies for video indexing, *Signal Processing: Image Communication*, **22(7)**, 692–704.
- [5] **Ayache, S. and Quénot, G.**, 2008. Video corpus annotation using active learning, *Advances in Information Retrieval*, Springer, pp.187–198.
- [6] **Joshi, H., Bayrak, C. and Xu, X.**, 2006. UALR at TREC: Blog Track., TREC.
- [7] **Bilgic, M., Mihalkova, L. and Getoor, L.**, 2010. Active learning for networked data, Proceedings of the 27th International Conference on Machine Learning (ICML-10), pp.79–86.
- [8] **Raghavan, H., Madani, O. and Jones, R.**, 2006. Active learning with feedback on features and instances, *The Journal of Machine Learning Research*, **7**, 1655–1686.
- [9] Collaborative Annotation of multi-Modal, multi-Lingual and multi-media documents, <http://www.chistera.eu/projects/camomile>.
- [10] **Poignant, J., Besacier, L., Quénot, G. and Thollard, F.**, 2012. From text detection in videos to person identification, Multimedia and Expo (ICME), 2012 IEEE International Conference on, IEEE, pp.854–859.
- [11] **Anthimopoulos, M., Gatos, B. and Pratikakis, I.**, 2010. A two-stage scheme for text detection in video images, *Image and Vision Computing*, **28(9)**, 1413–1426.
- [12] **Chen, S. and Gopalakrishnan, P.**, 1998. Speaker, environment and channel change detection and clustering via the Bayesian Information Criterion, Proc. DARPA Broadcast News Transcription and Understanding Workshop, Virginia, USA, p. 8.

- [13] **Guillaumin, M., Mensink, T., Verbeek, J. and Schmid, C.**, 2012. Face recognition from caption-based supervision, *International Journal of Computer Vision*, **96(1)**, 64–82.
- [14] **Settles, B.**, 2012. Active learning, *Synthesis Lectures on Artificial Intelligence and Machine Learning*, **6(1)**, 1–114.
- [15] **Valiant, L.G.**, 1984. A theory of the learnable, *Communications of the ACM*, **27(11)**, 1134–1142.
- [16] **Angluin, D.**, 1988. Queries and concept learning, *Machine learning*, **2(4)**, 319–342.
- [17] **Settles, B.**, 2010. Active learning literature survey, *University of Wisconsin, Madison*, **52**, 55–66.
- [18] **Lewis, D.D. and Catlett, J.**, 1994. Heterogenous Uncertainty Sampling for Supervised Learning., ICML, volume 94, pp.148–156.
- [19] **Shannon, C.E.**, 2001. A mathematical theory of communication, *ACM SIGMOBILE Mobile Computing and Communications Review*, **5(1)**, 3–55.
- [20] **Cohn, D., Atlas, L. and Ladner, R.**, 1994. Improving generalization with active learning, *Machine learning*, **15(2)**, 201–221.
- [21] **Seung, H.S., Opper, M. and Sompolinsky, H.**, 1992. Query by committee, Proceedings of the fifth annual workshop on Computational learning theory, ACM, pp.287–294.
- [22] **Breiman, L.**, 1996. Bagging predictors, *Machine learning*, **24(2)**, 123–140.
- [23] **Kullback, S. and Leibler, R.A.**, 1951. On information and sufficiency, *The Annals of Mathematical Statistics*, 79–86.
- [24] **Meilă, M.**, 2003. Comparing clusterings by the variation of information, Learning theory and kernel machines, Springer, pp.173–187.
- [25] **Wallace, D.L.**, 1983. Comment, *Journal of the American Statistical Association*, **78(383)**, 569–576.
- [26] **Rand, W.M.**, 1971. Objective criteria for the evaluation of clustering methods, *Journal of the American Statistical association*, **66(336)**, 846–850.
- [27] **Mirkin, B.**, 1998. Mathematical classification and clustering: From how to what and why, Springer.
- [28] **Meilă, M. and Heckerman, D.**, 2001. An experimental comparison of model-based clustering methods, *Machine learning*, **42(1-2)**, 9–29.
- [29] **Dongen, S.**, 2000. Performance criteria for graph clustering and Markov cluster experiments.

- [30] **Amigó, E., Gonzalo, J., Artiles, J. and Verdejo, F.**, 2009. A comparison of extrinsic clustering evaluation metrics based on formal constraints, *Information retrieval*, **12(4)**, 461–486.
- [31] **Iwama, K. and Miyazaki, S.**, 2008. A survey of the stable marriage problem and its variants, Informatics Education and Research for Knowledge-Circulating Society, 2008. ICKS 2008. International Conference on, IEEE, pp.131–136.
- [32] **Krumpelman, C. and Ghosh, J.**, 2007. Matching and visualization of multiple overlapping clusterings of microarray data, Computational Intelligence and Bioinformatics and Computational Biology, 2007. CIBCB'07. IEEE Symposium on, IET, pp.121–126.
- [33] **Cohn, D.A., Ghahramani, Z. and Jordan, M.I.**, 1996. Active learning with statistical models, *arXiv preprint cs/9603104*.
- [34] **Atlas, L.E., Cohn, D.A. and Ladner, R.E.**, 1990. Training connectionist networks with queries and selective sampling, Advances in neural information processing systems, pp.566–573.
- [35] **Settles, B. and Craven, M.**, 2008. An analysis of active learning strategies for sequence labeling tasks, Proceedings of the Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, pp.1070–1079.
- [36] **Xu, Z., Akella, R. and Zhang, Y.**, 2007. Incorporating diversity and density in active learning for relevance feedback, Springer.



## **APPENDICES**

### **APPENDIX A.1 : Video Based Results and Methods Ranking Heat Maps**



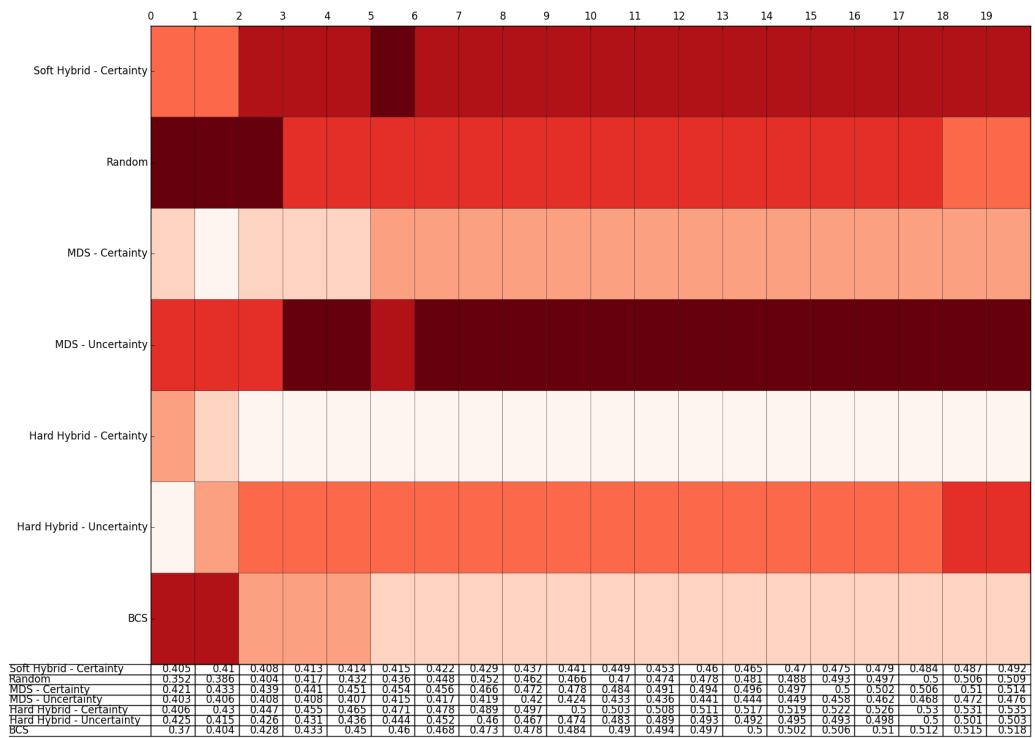


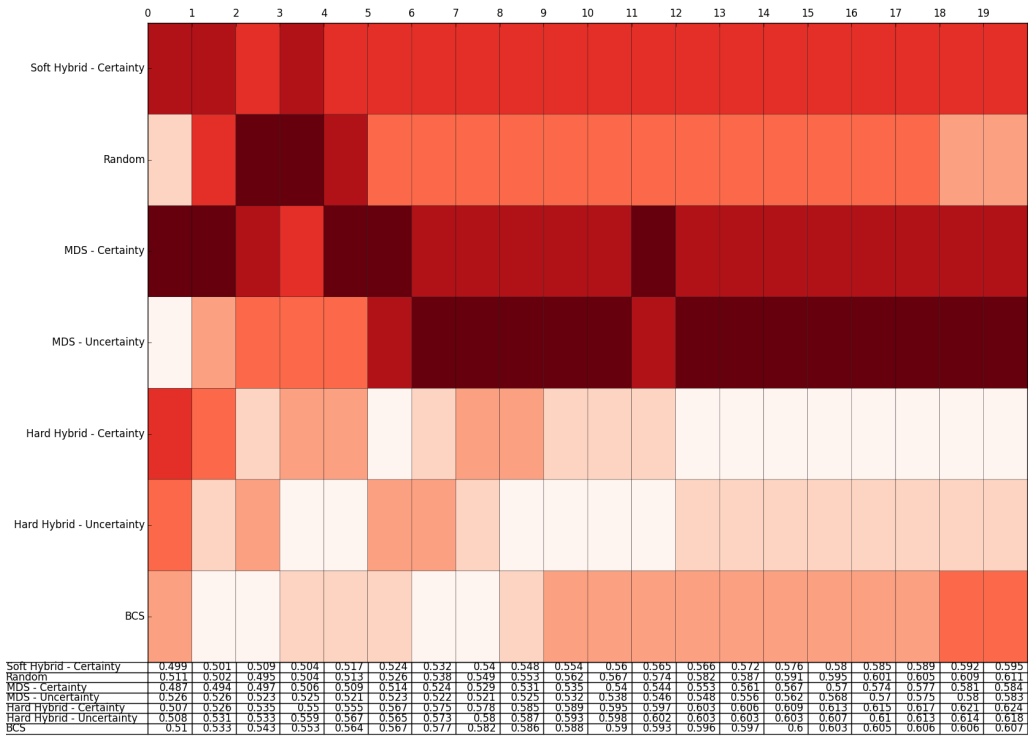
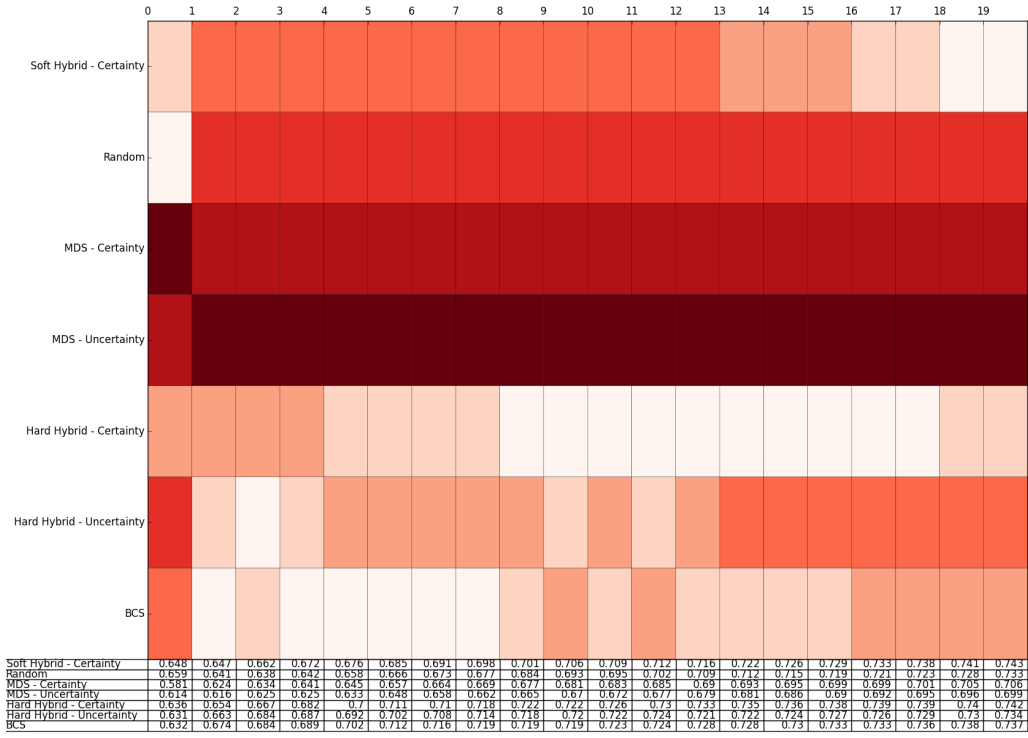
## APPENDIX A.1

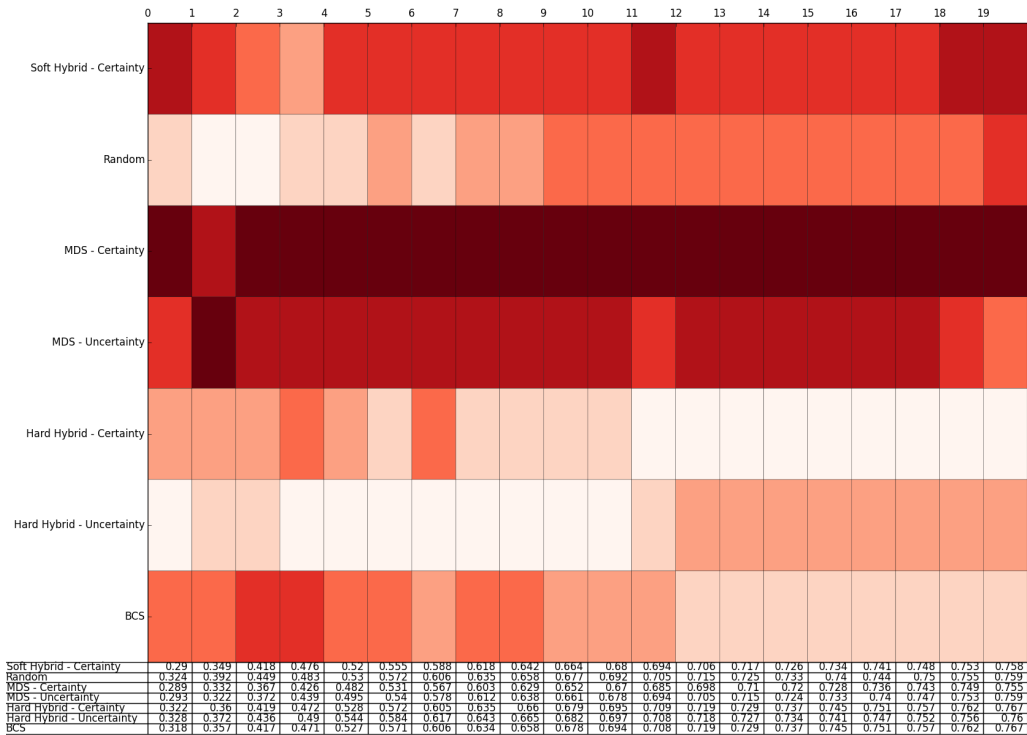
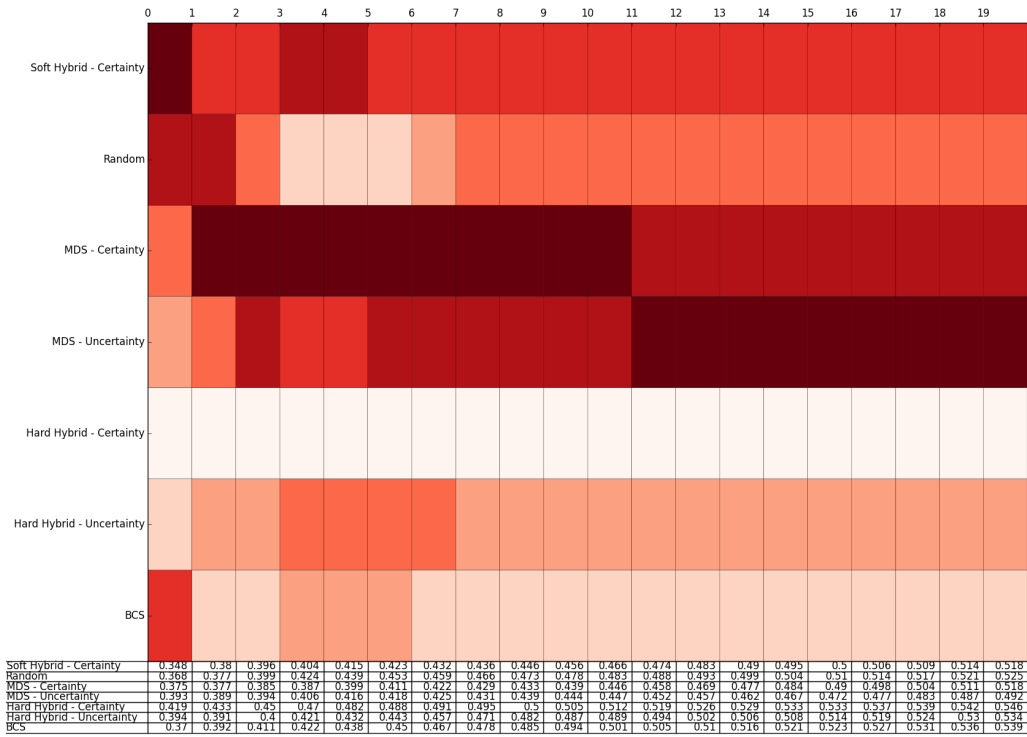
F-measure values and algorithm heat maps of videos.

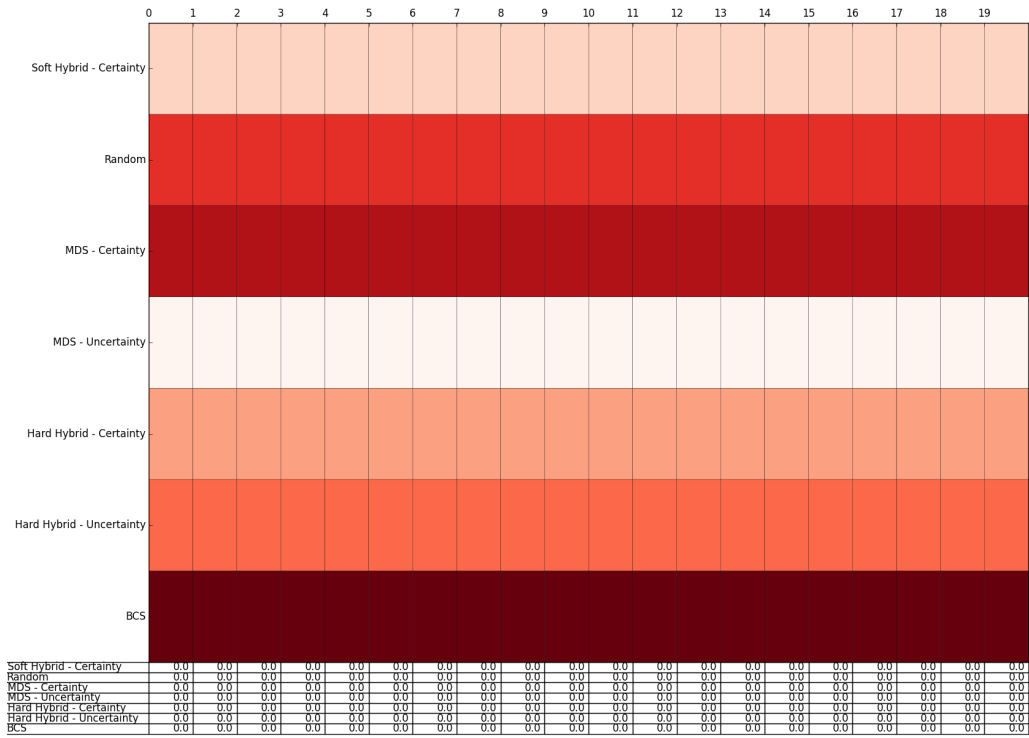
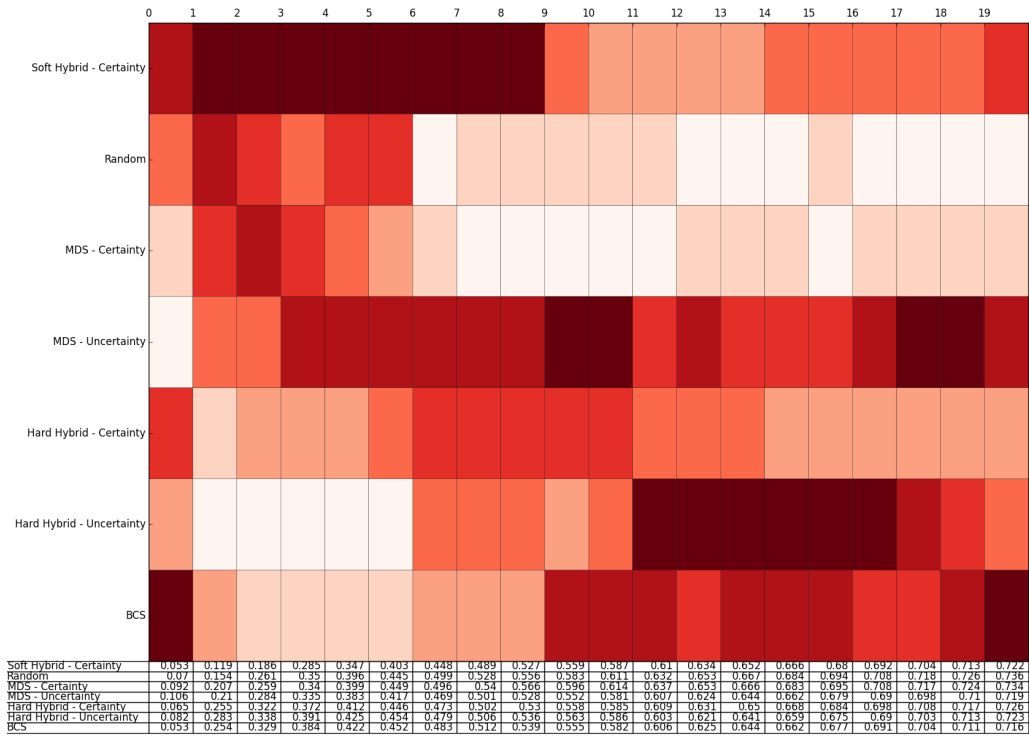
Face score using face annotation is given in red heat maps.

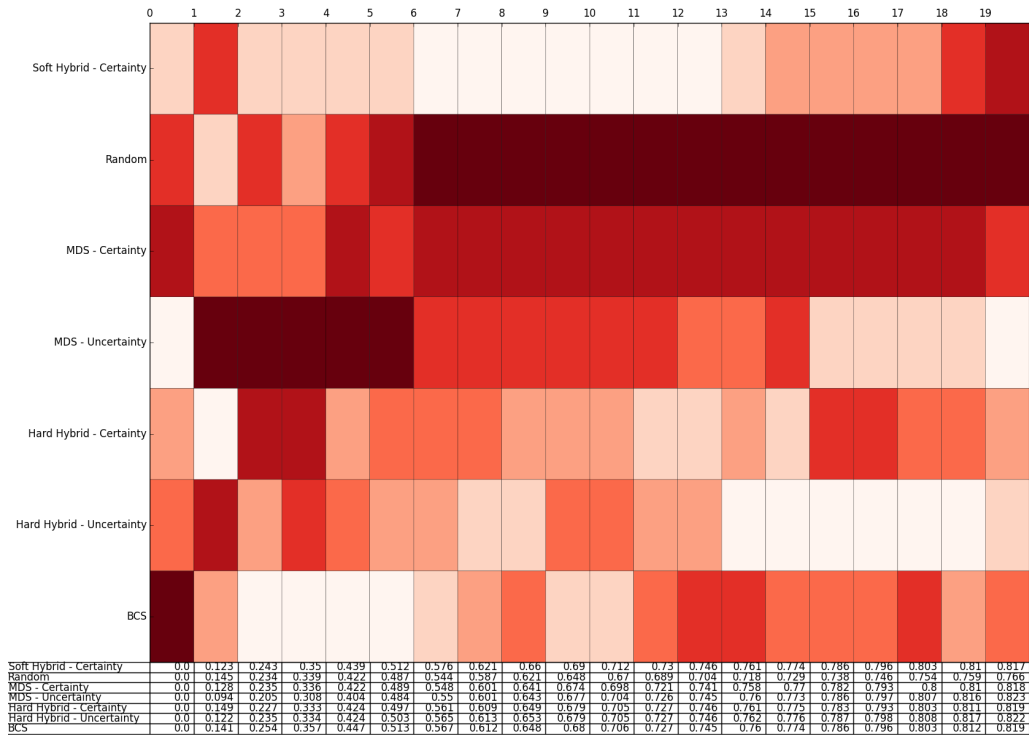
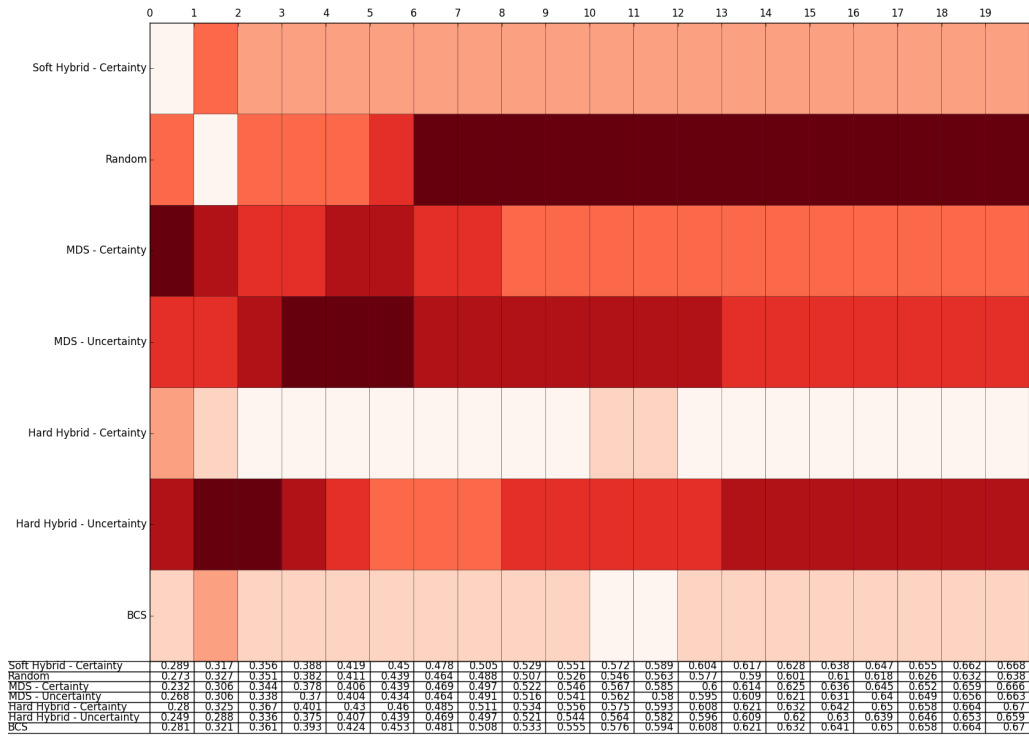
BFMTV\_BFMStory\_2012-01-10\_175800, Head Track: 759, Speaker Track: 372, OCR Name: 36 HvH 401X401 SvS 291X291 SvH 737

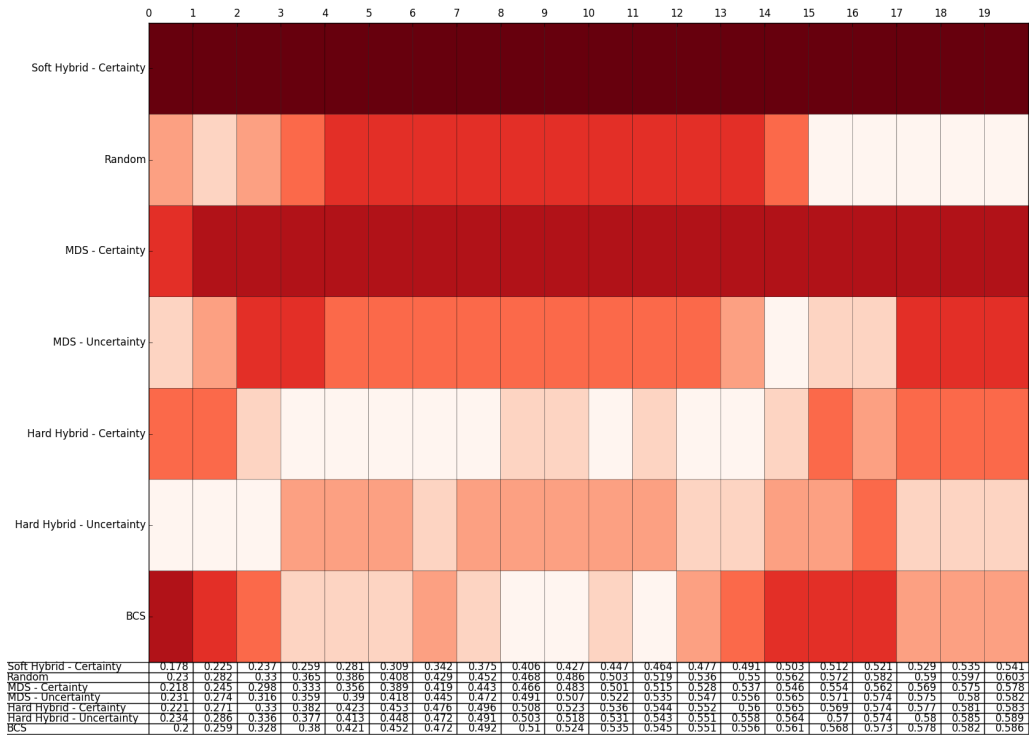
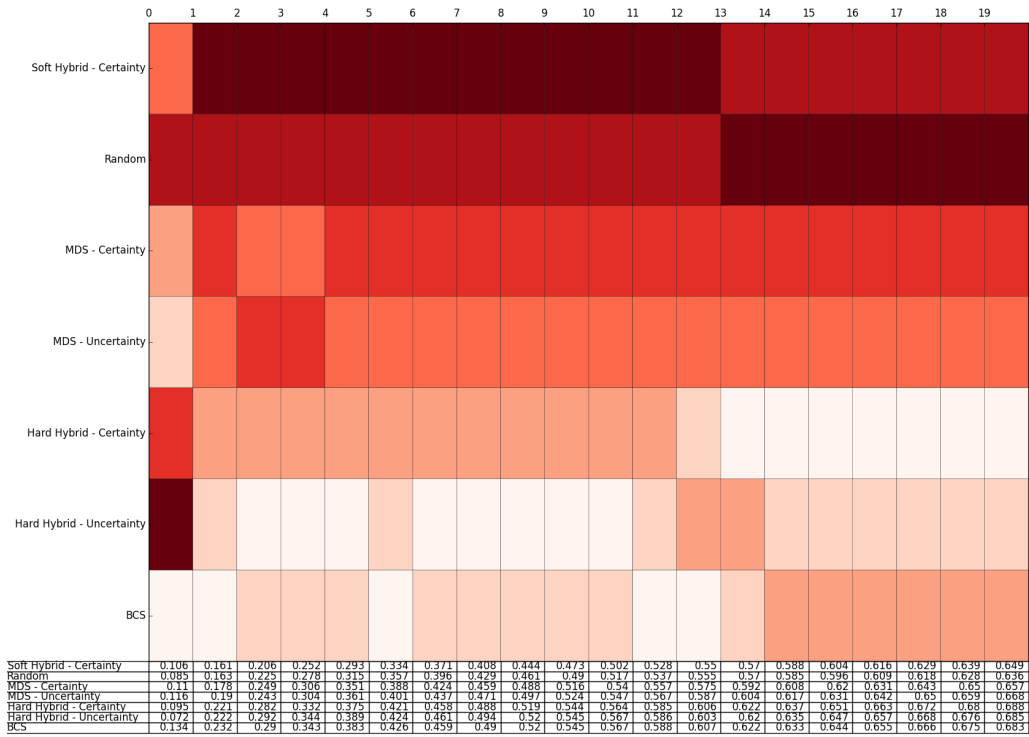


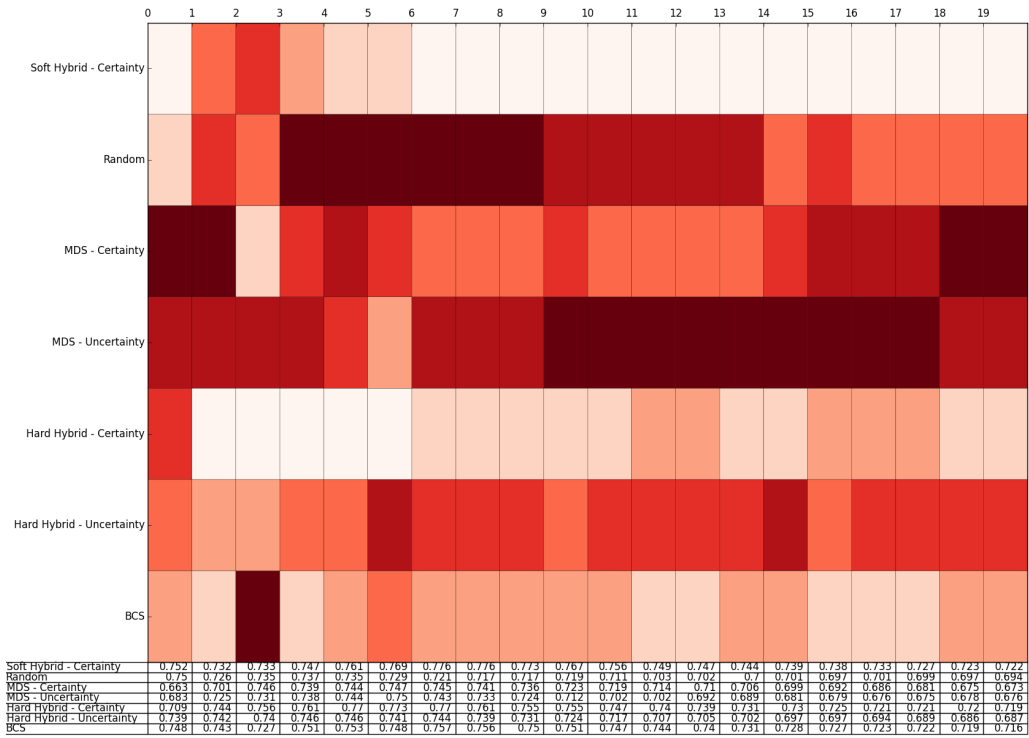
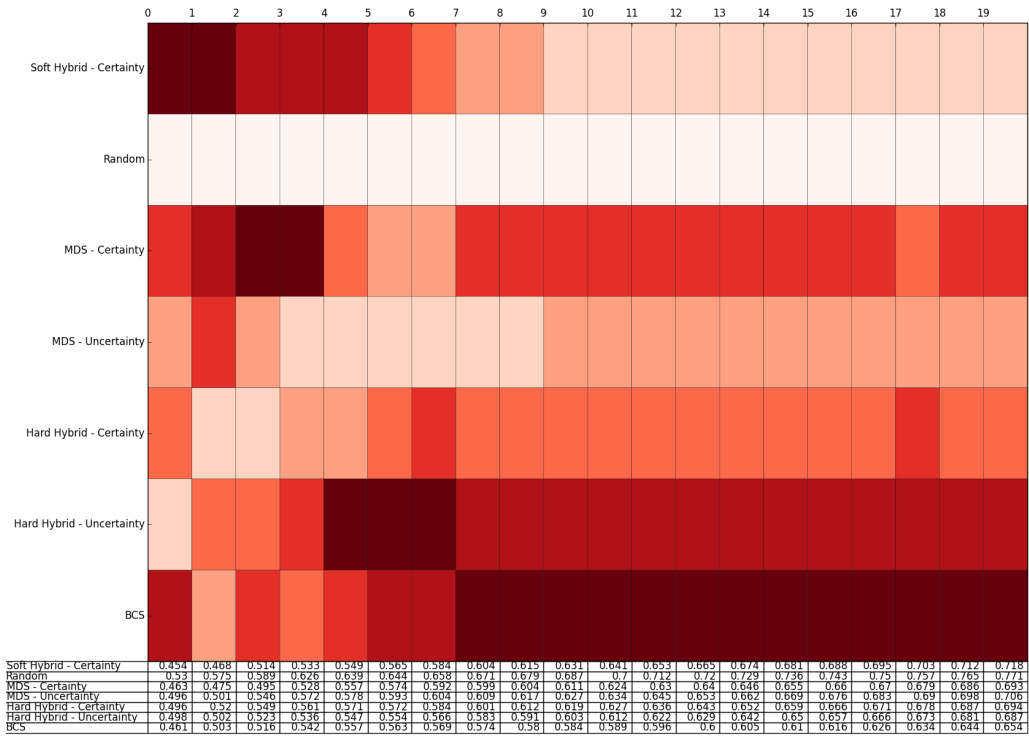


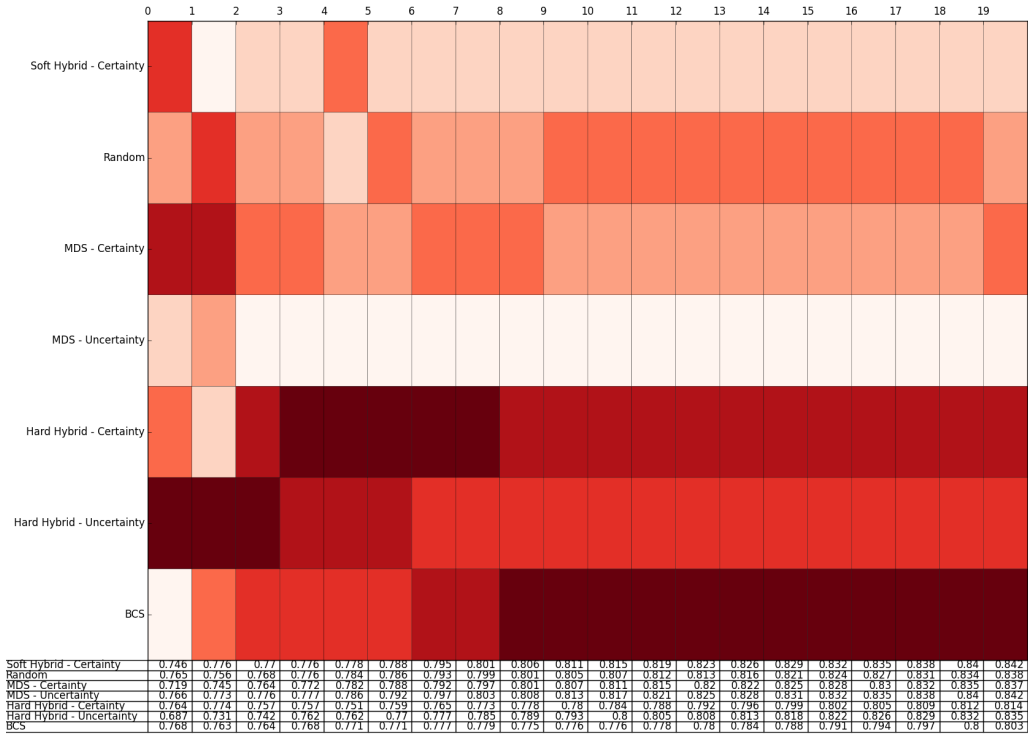
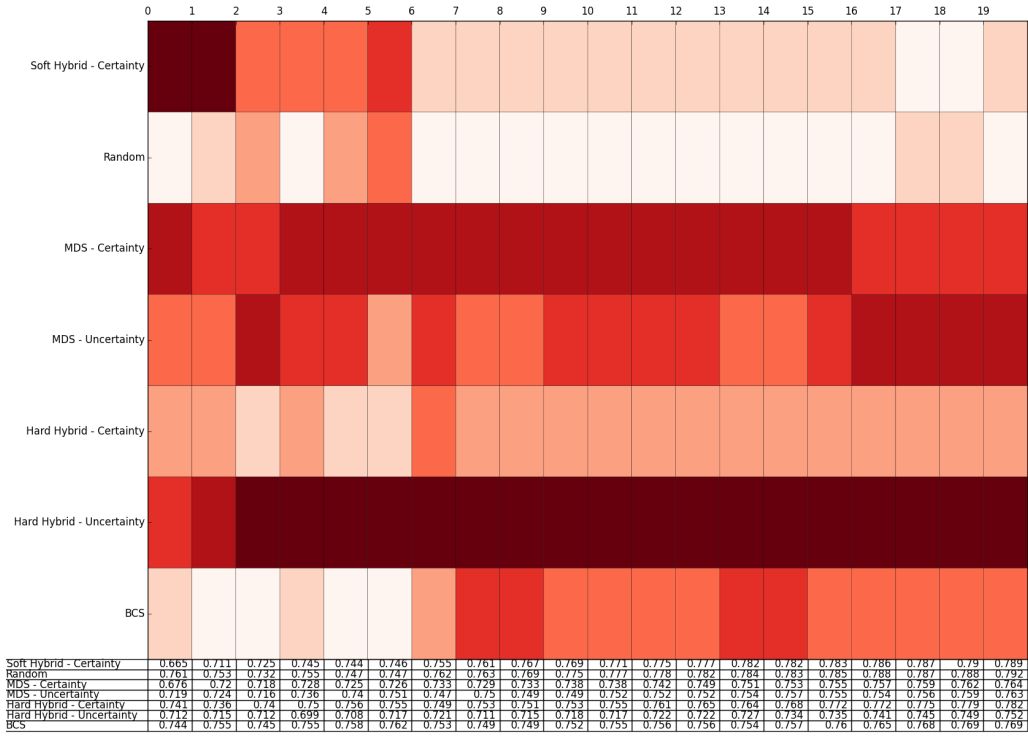




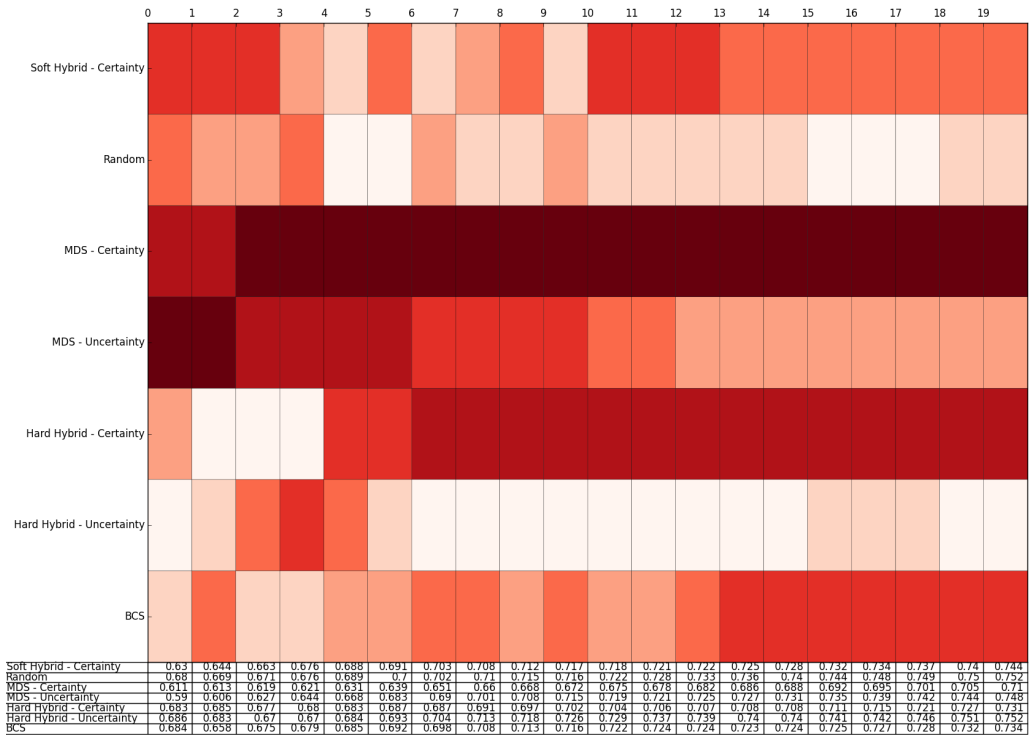
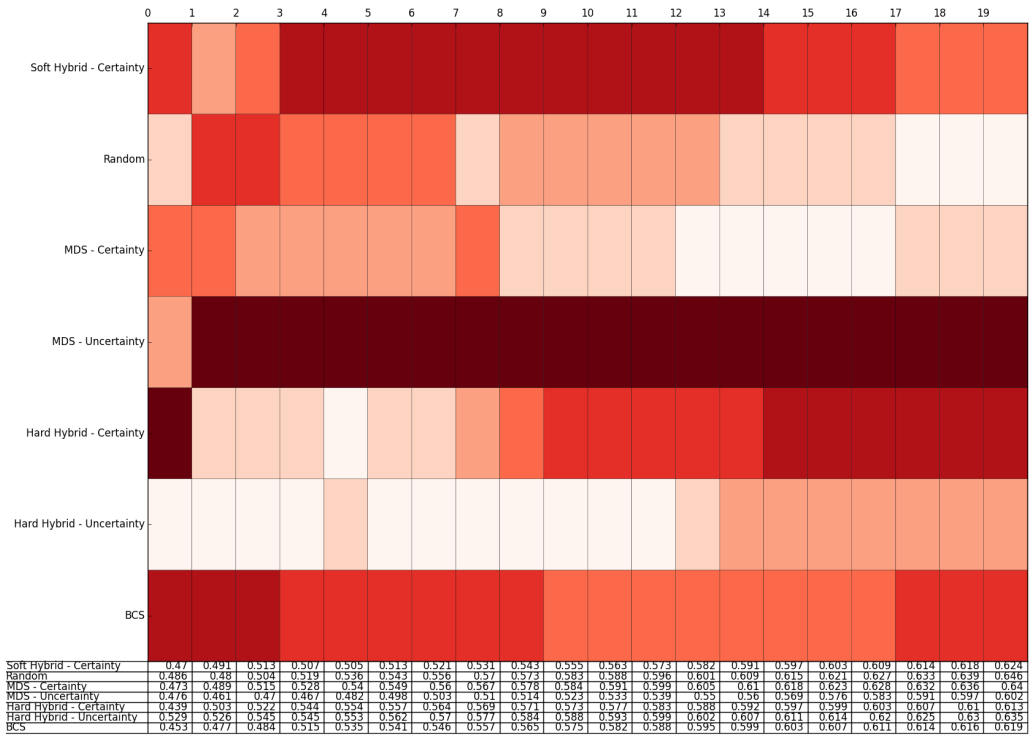


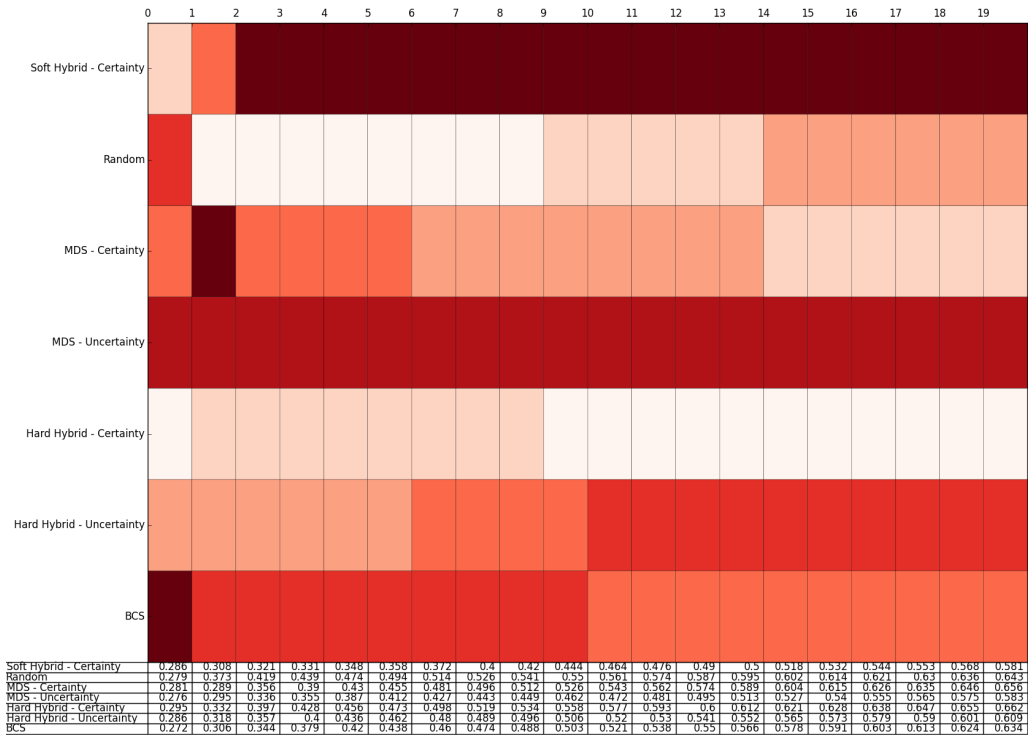
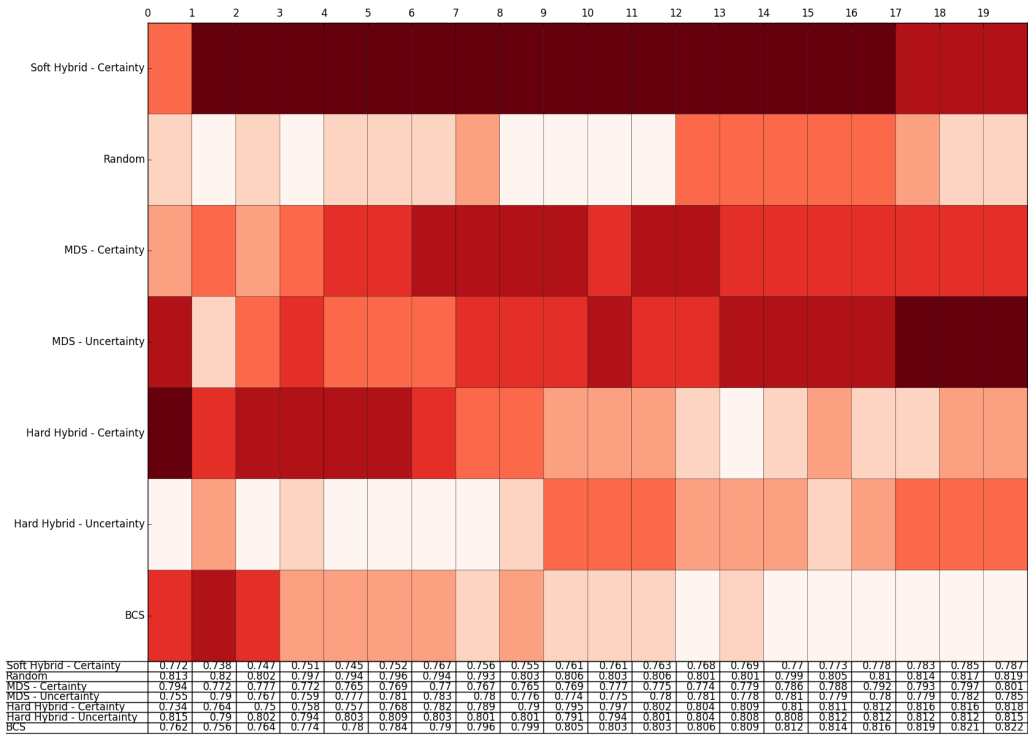


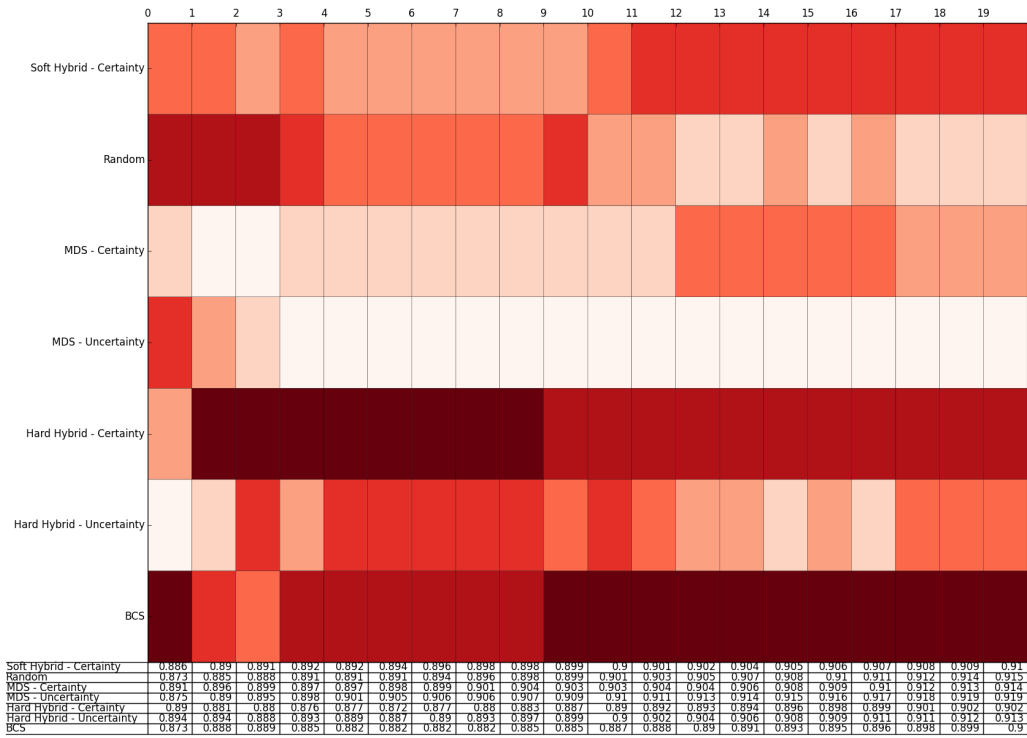
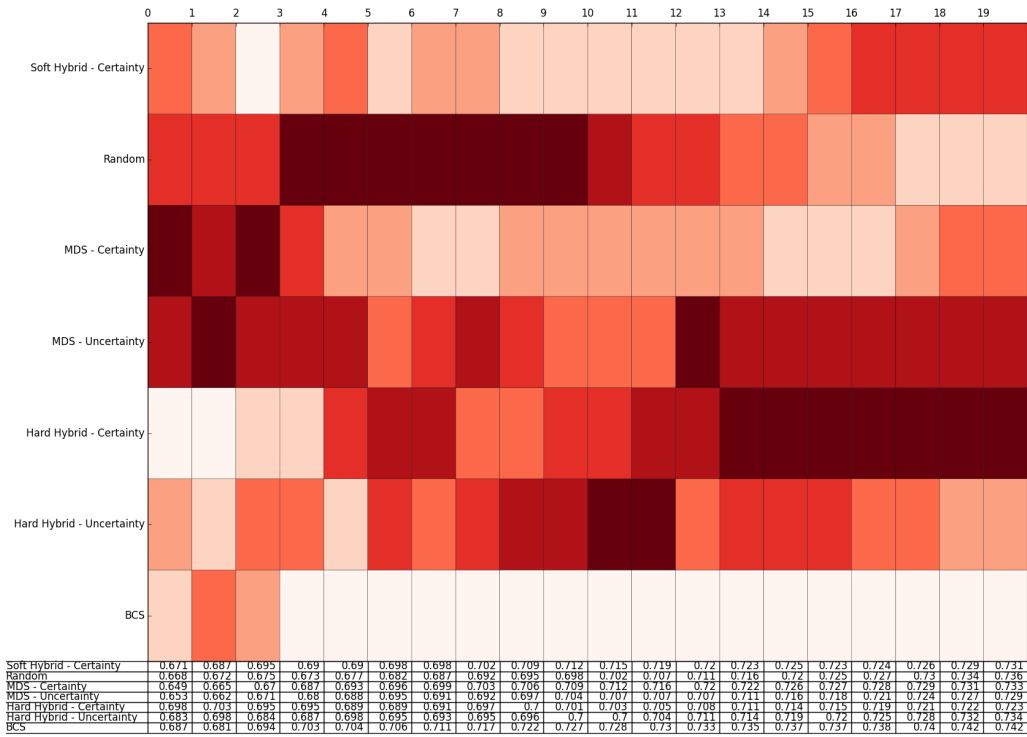


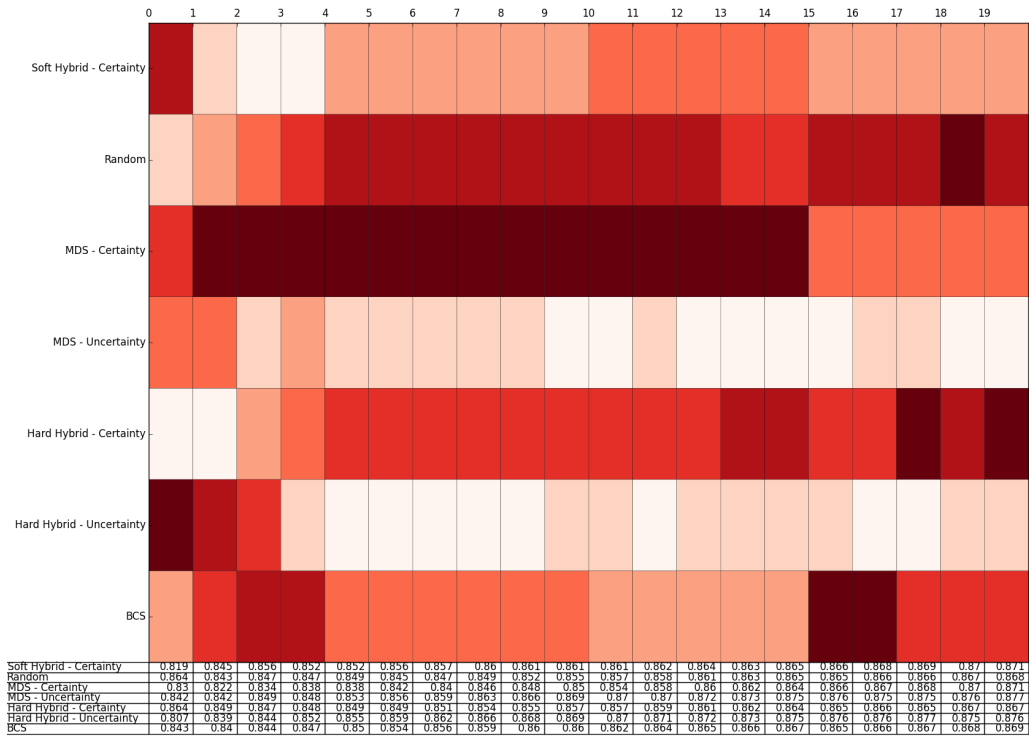
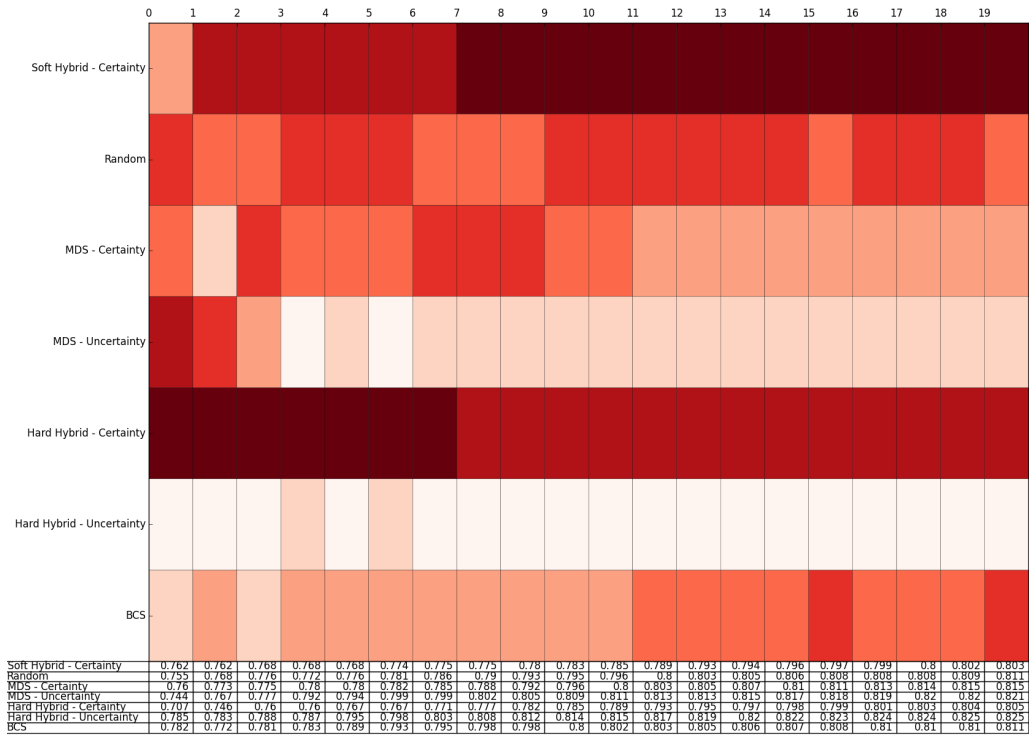


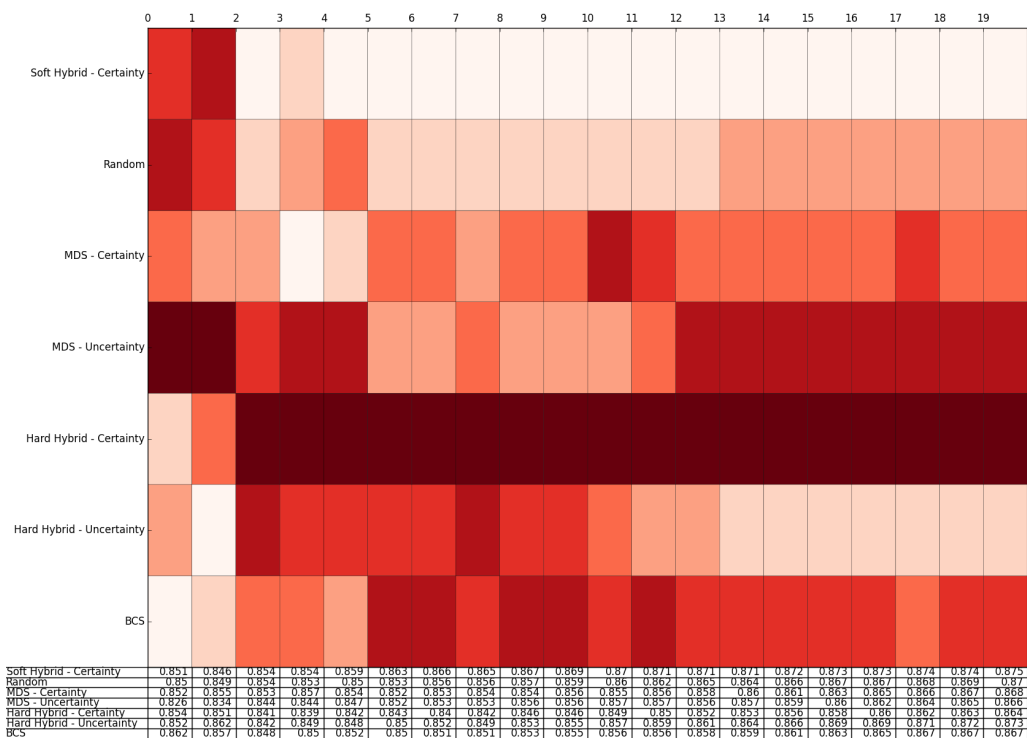
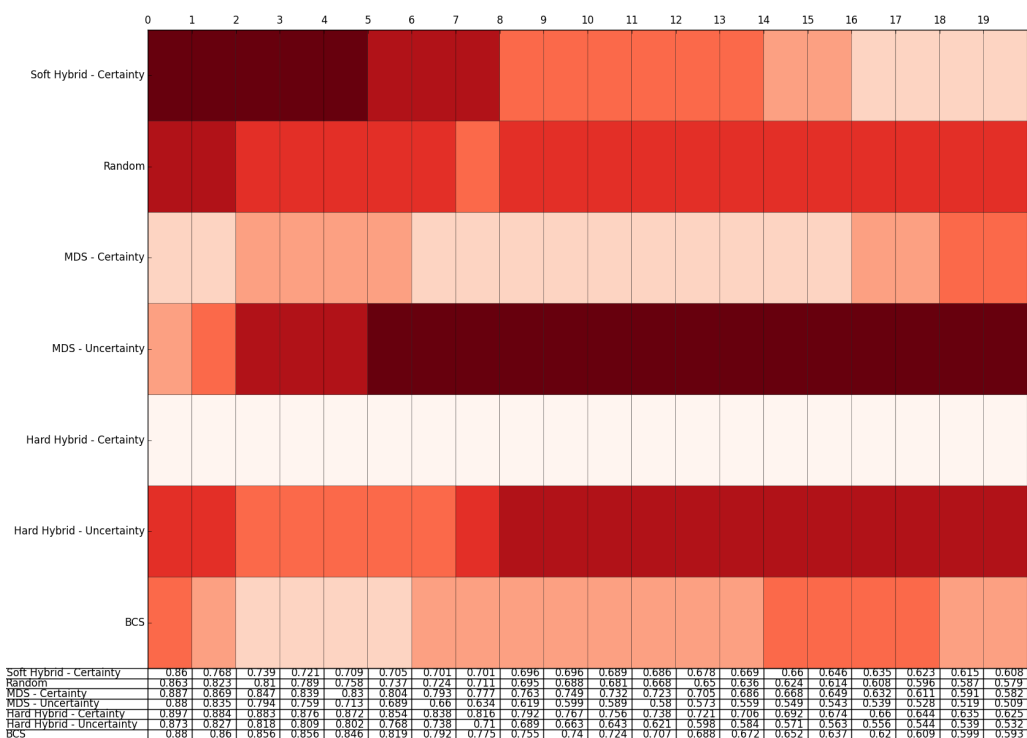


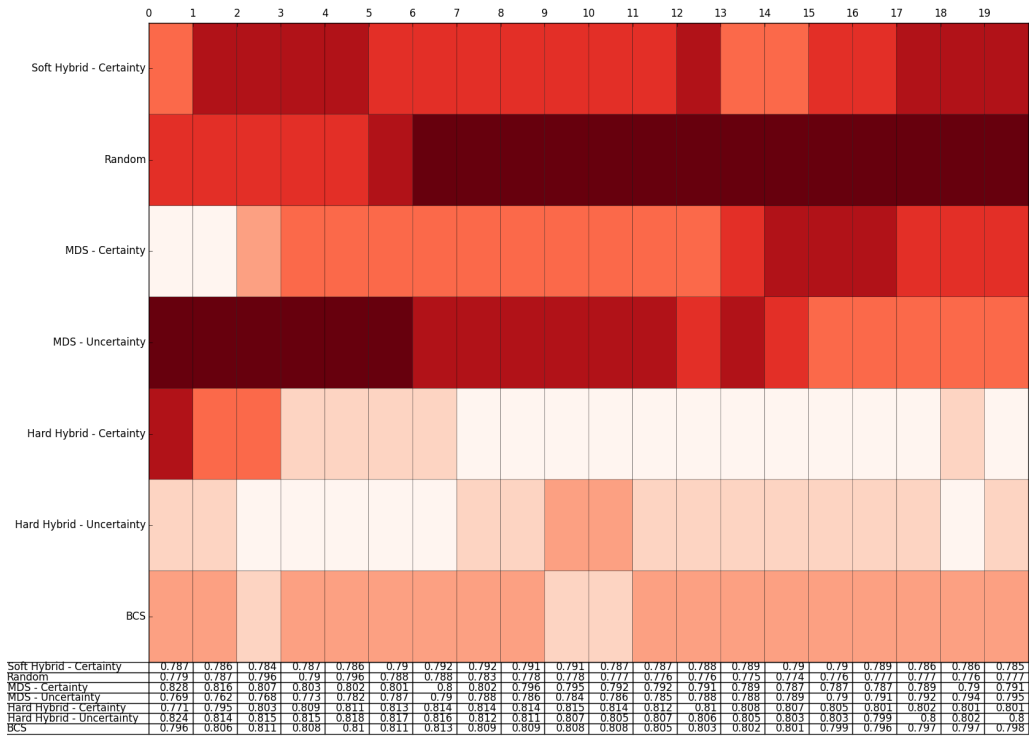
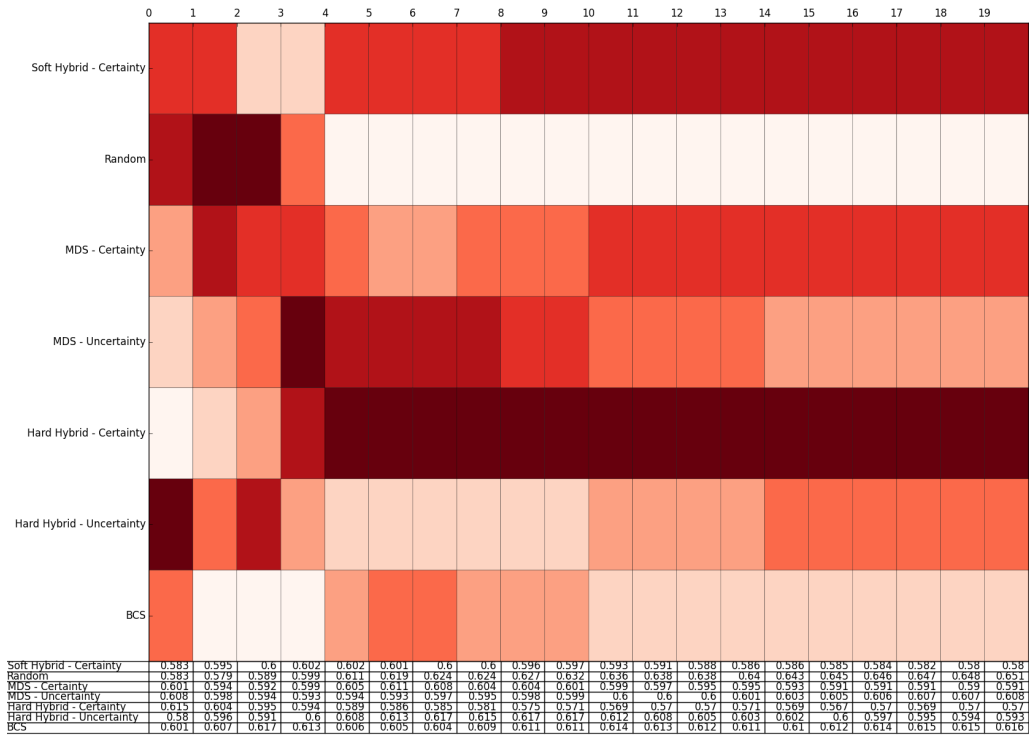


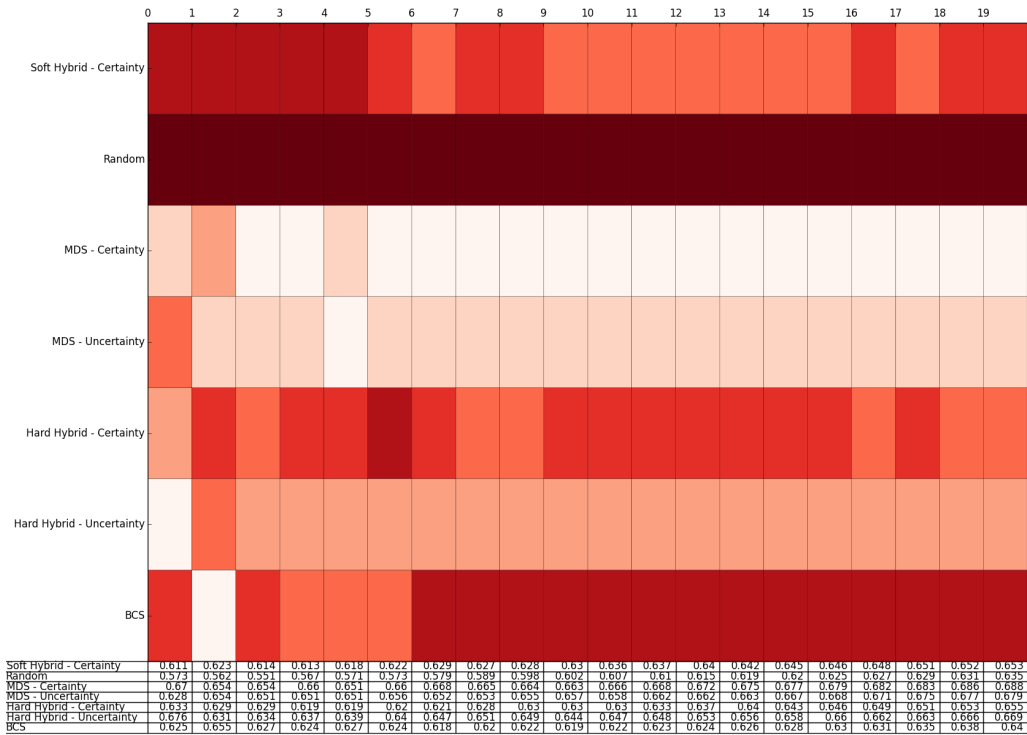




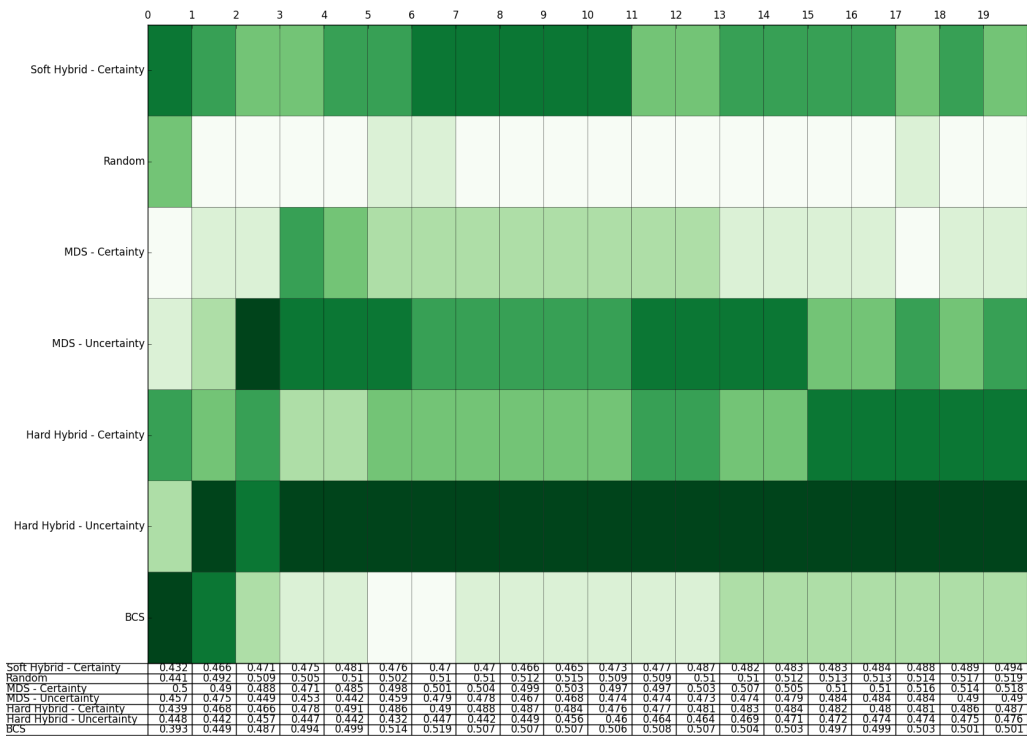


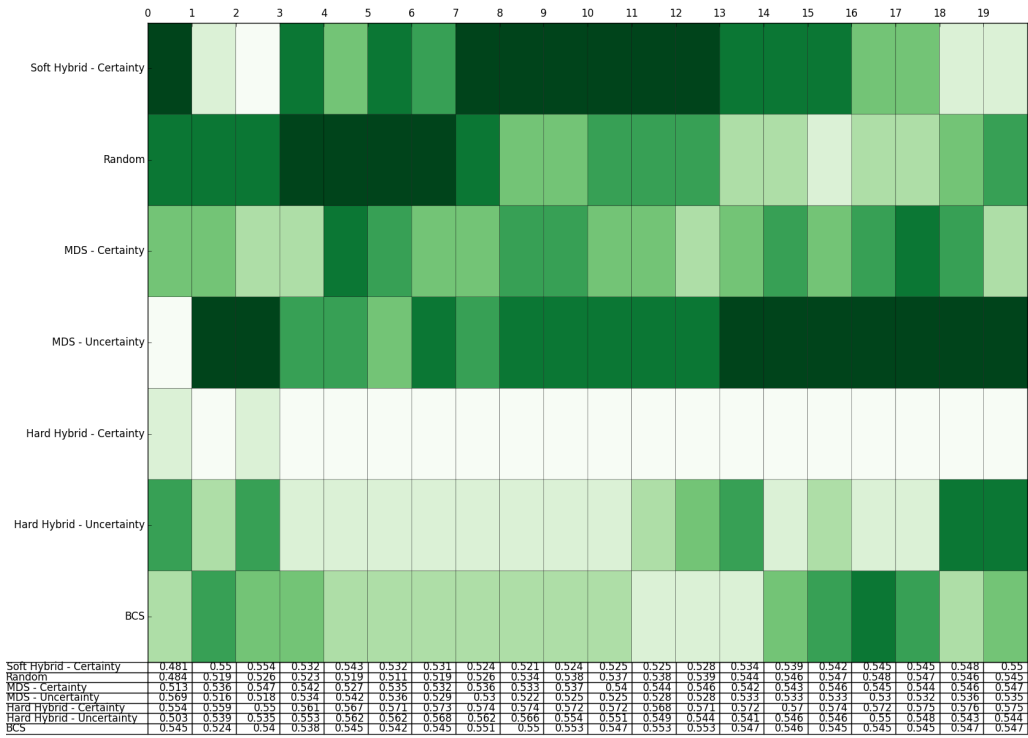
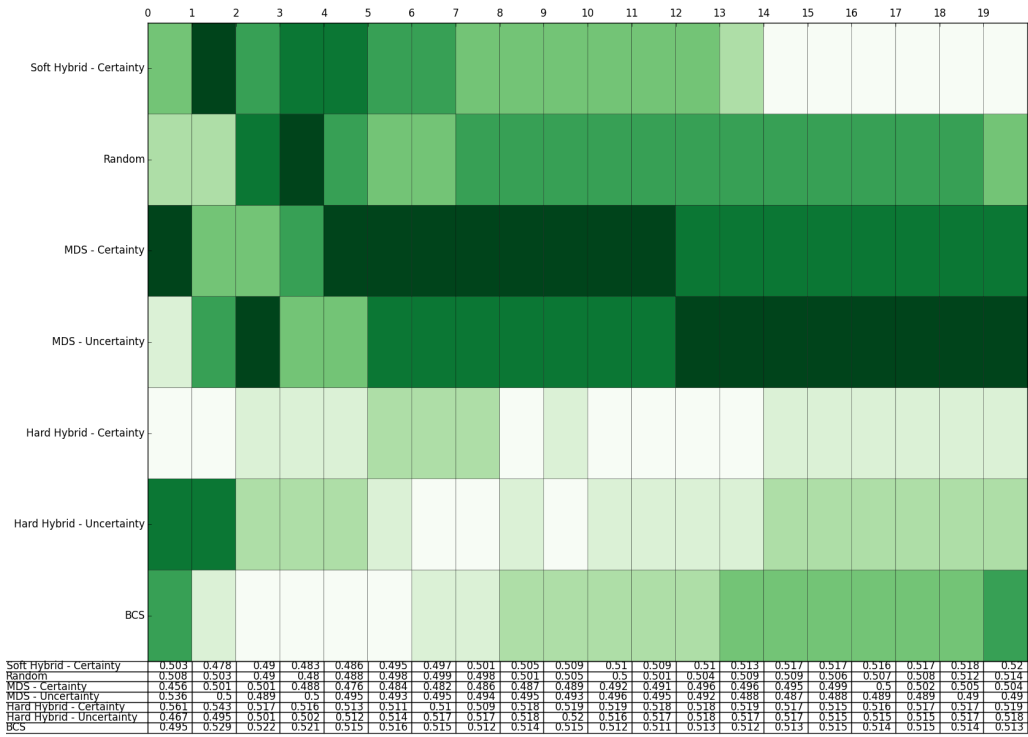




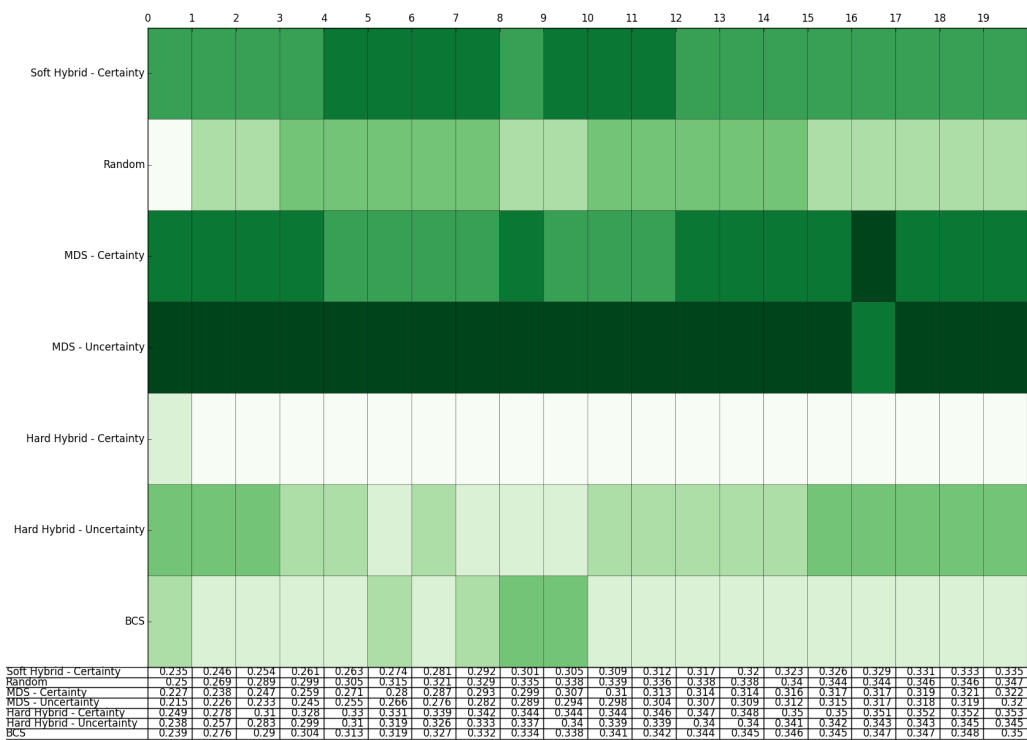
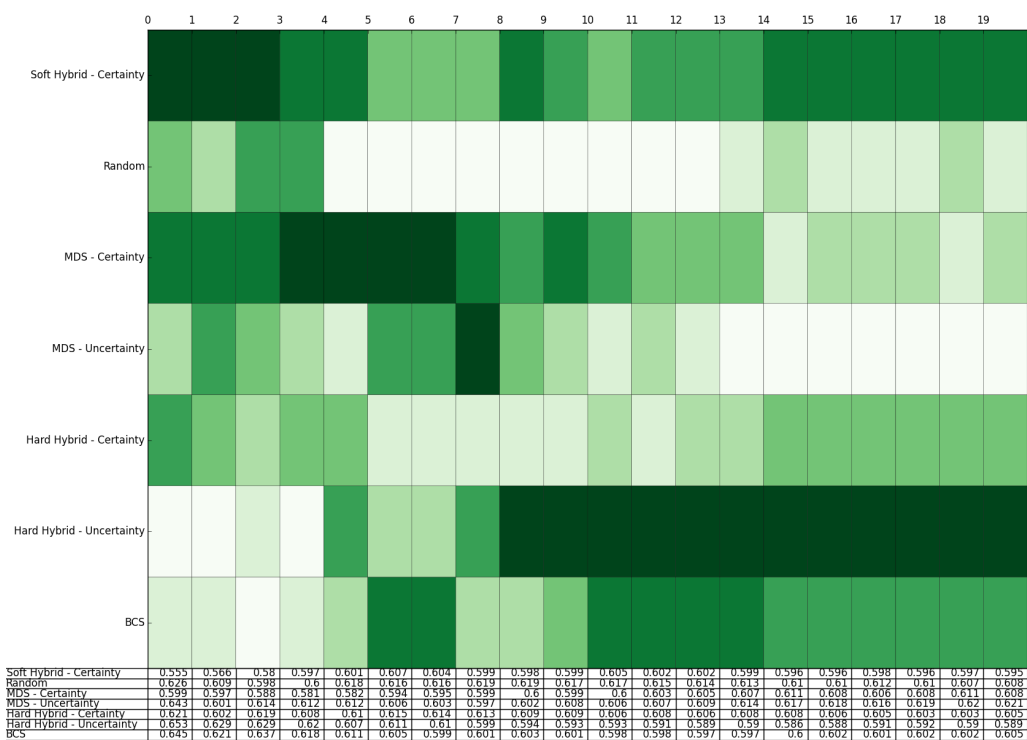


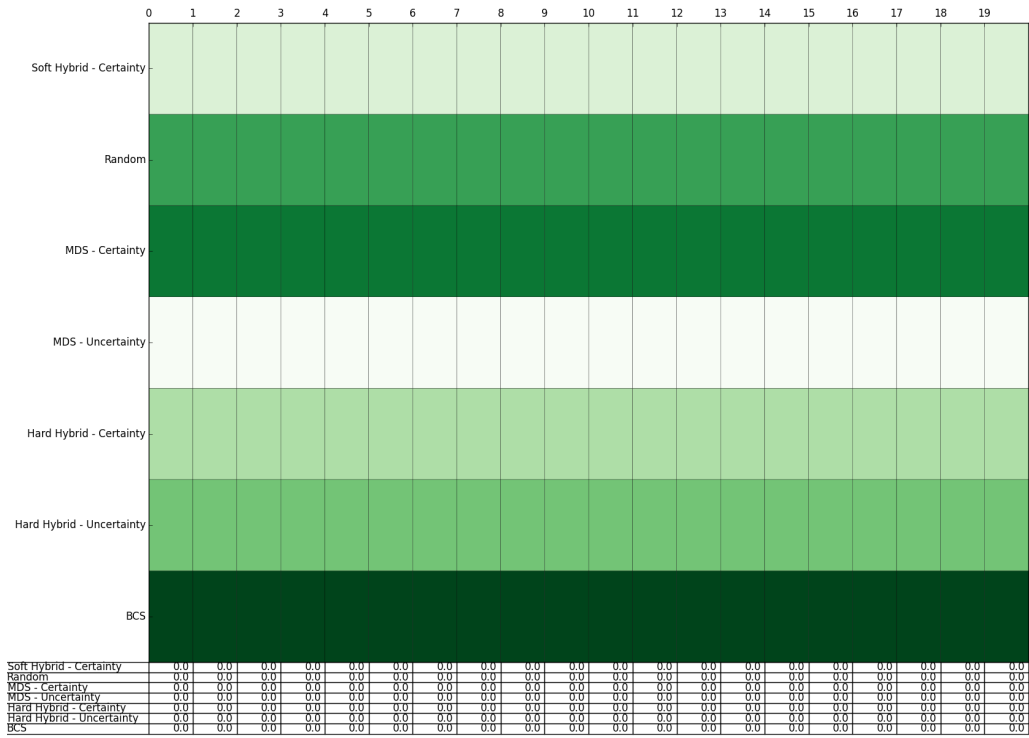
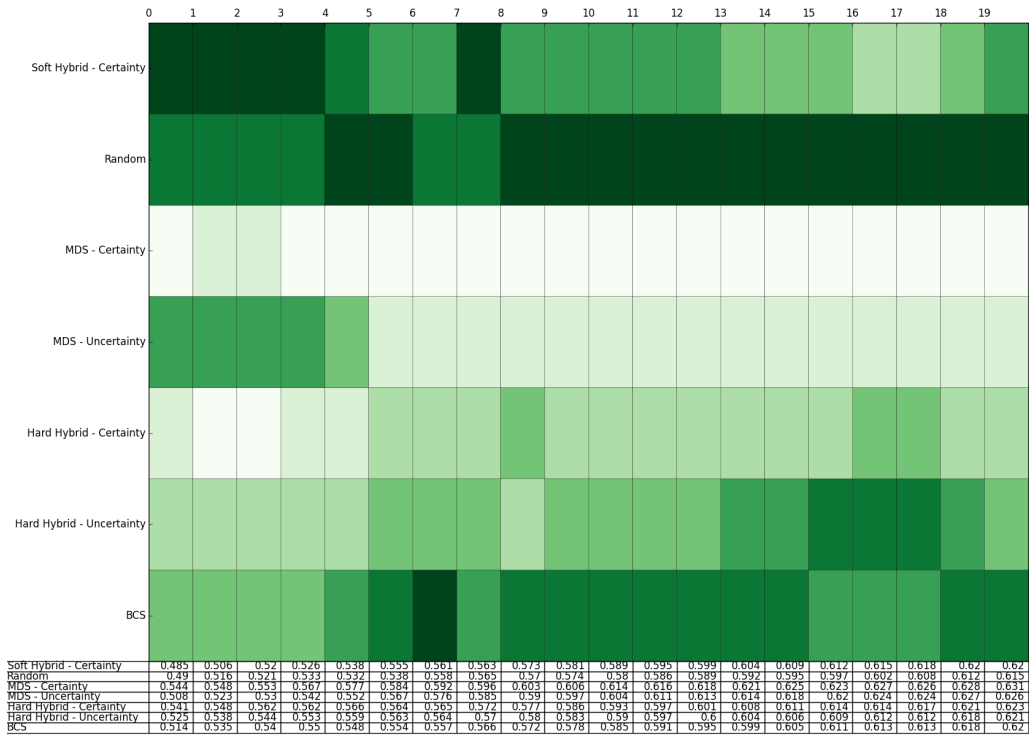
Speaker score using face annotation is given green in heat maps.

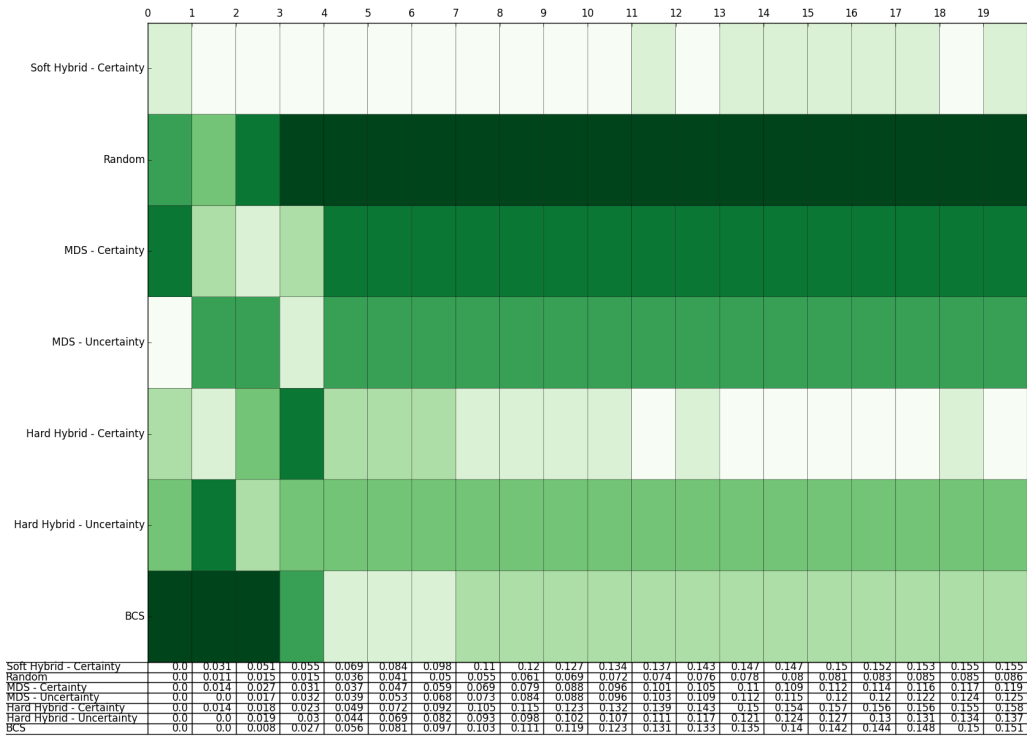
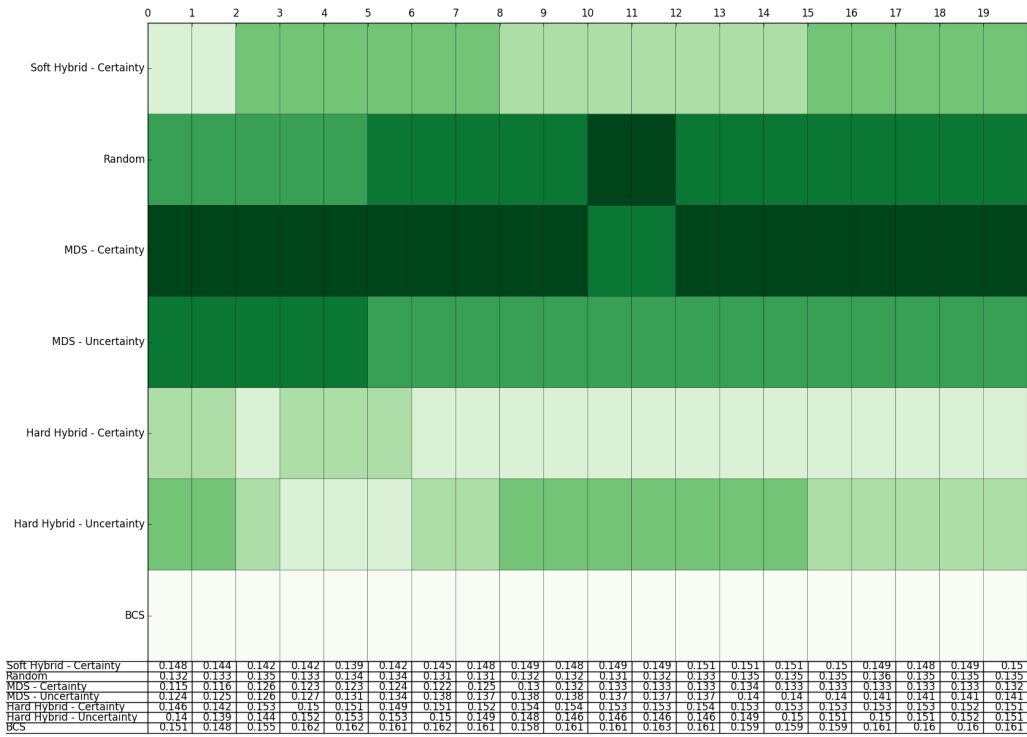


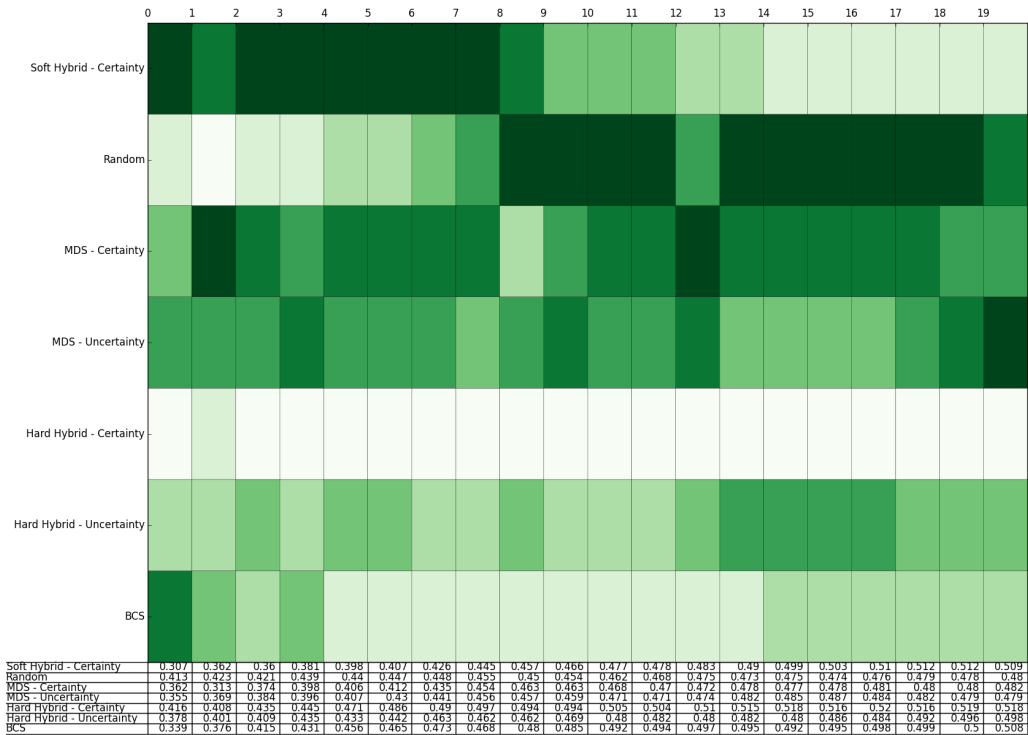
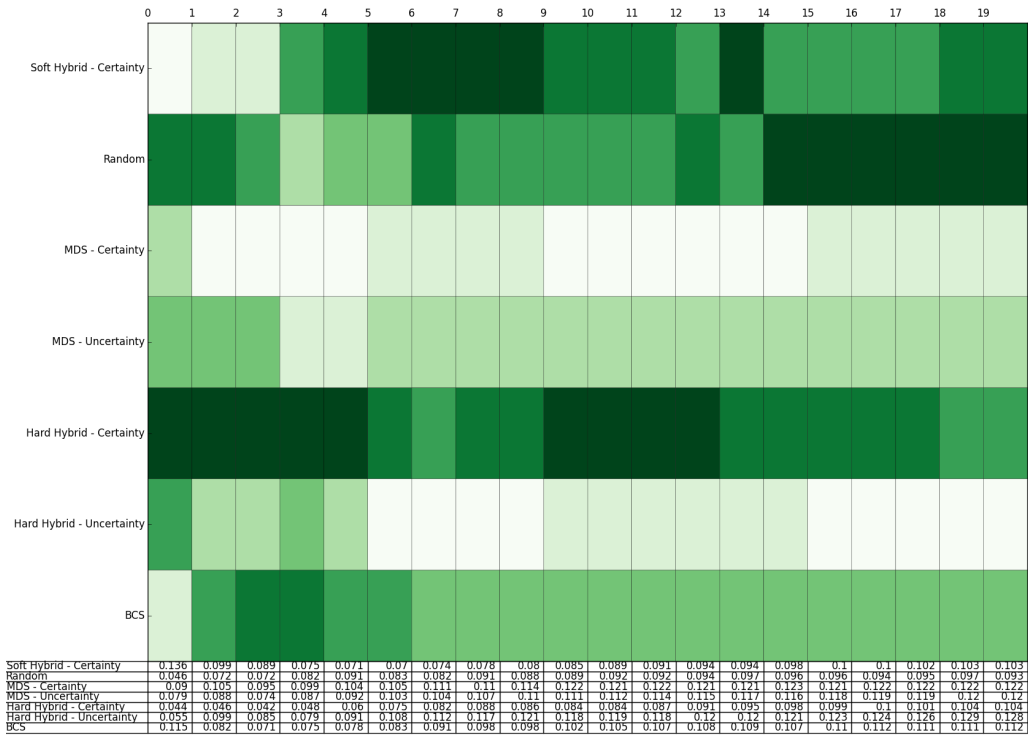


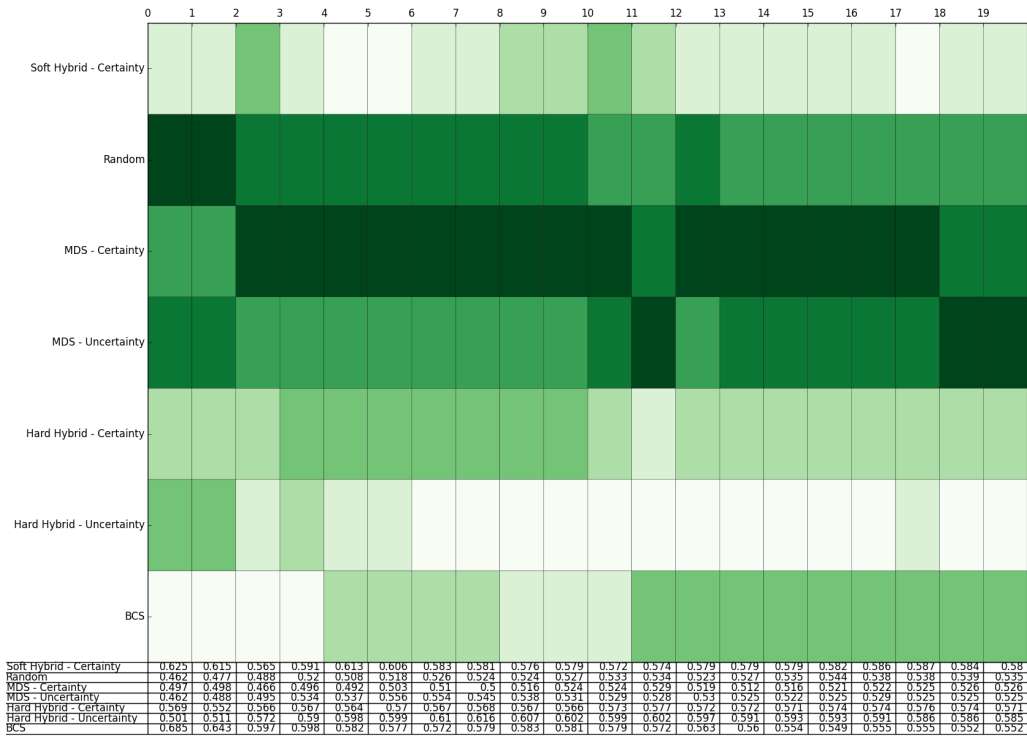
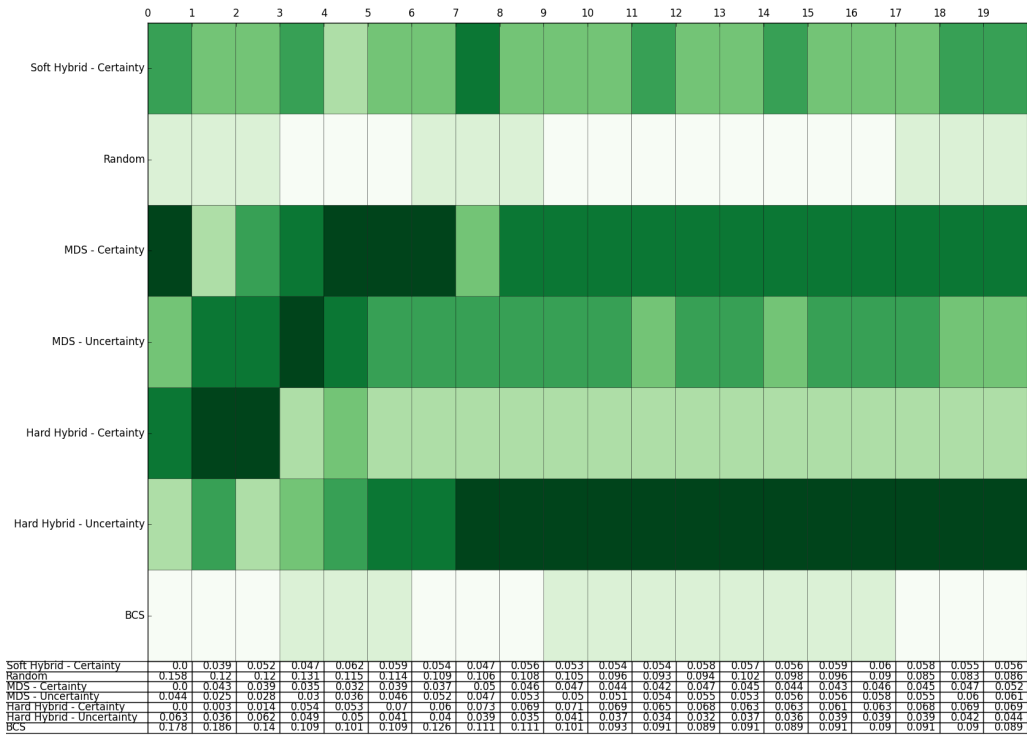


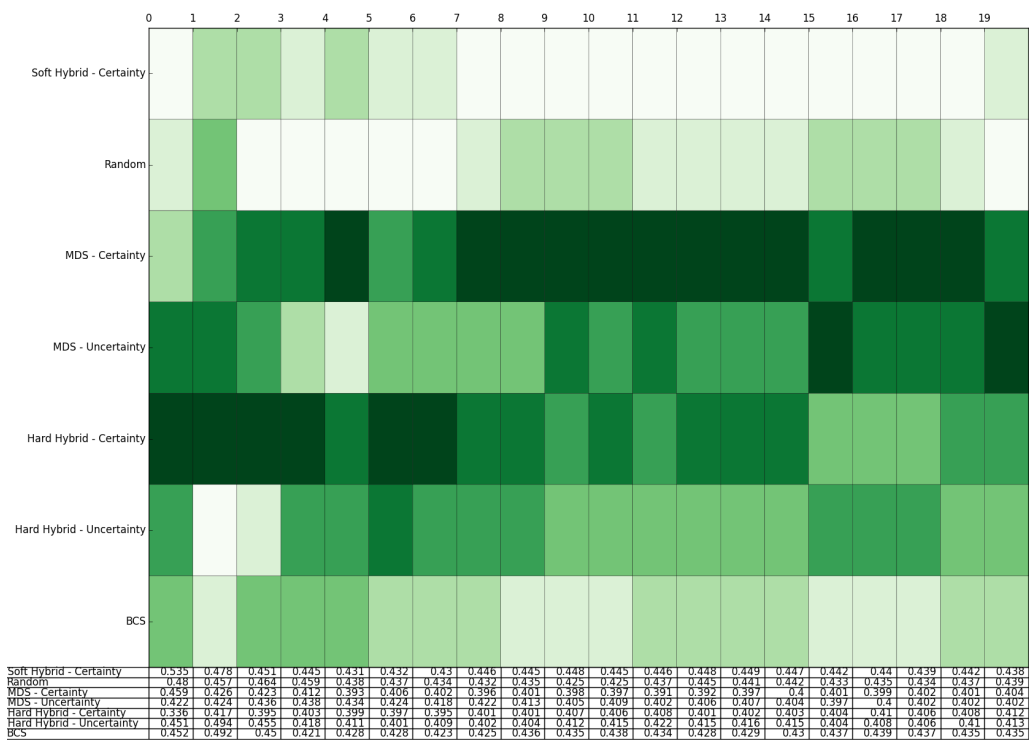
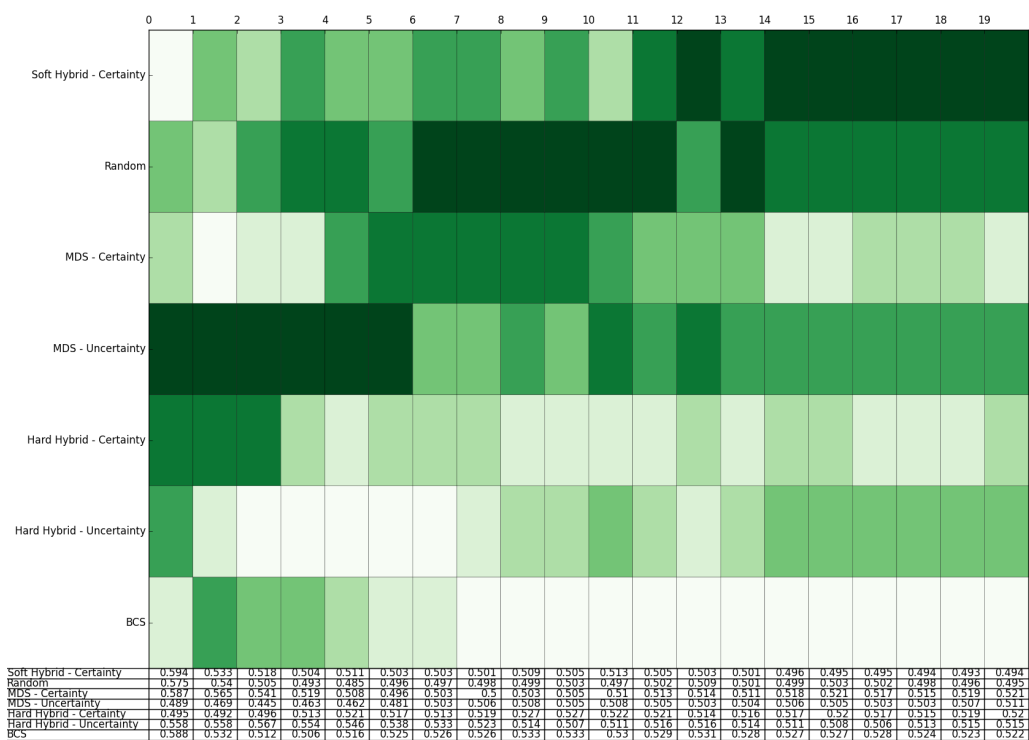


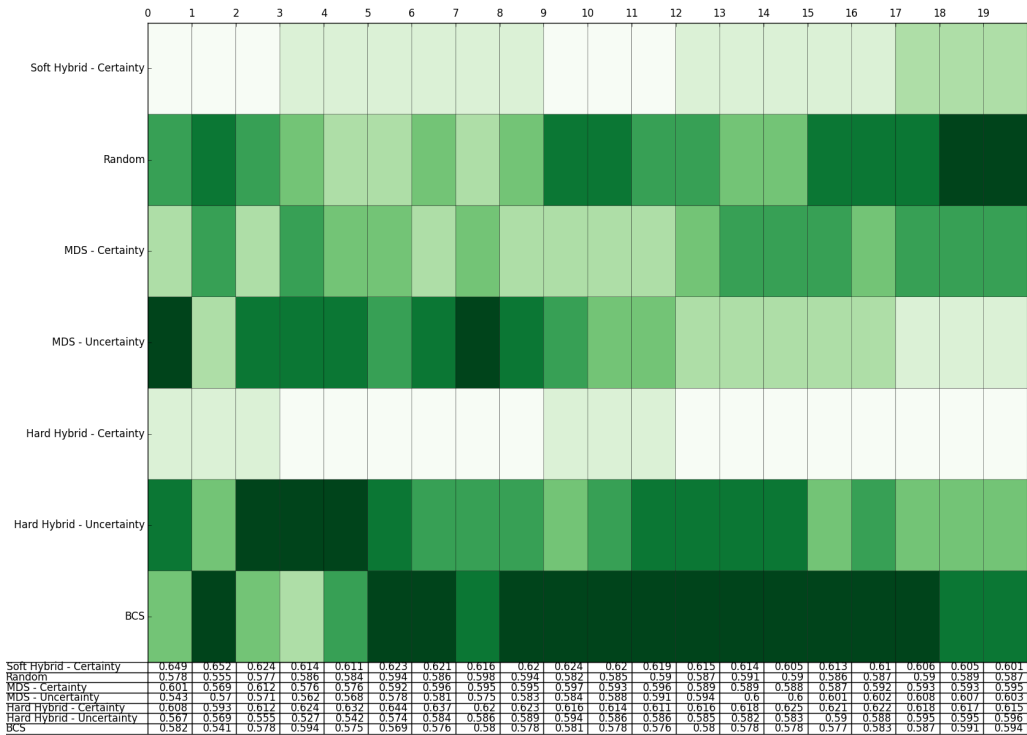
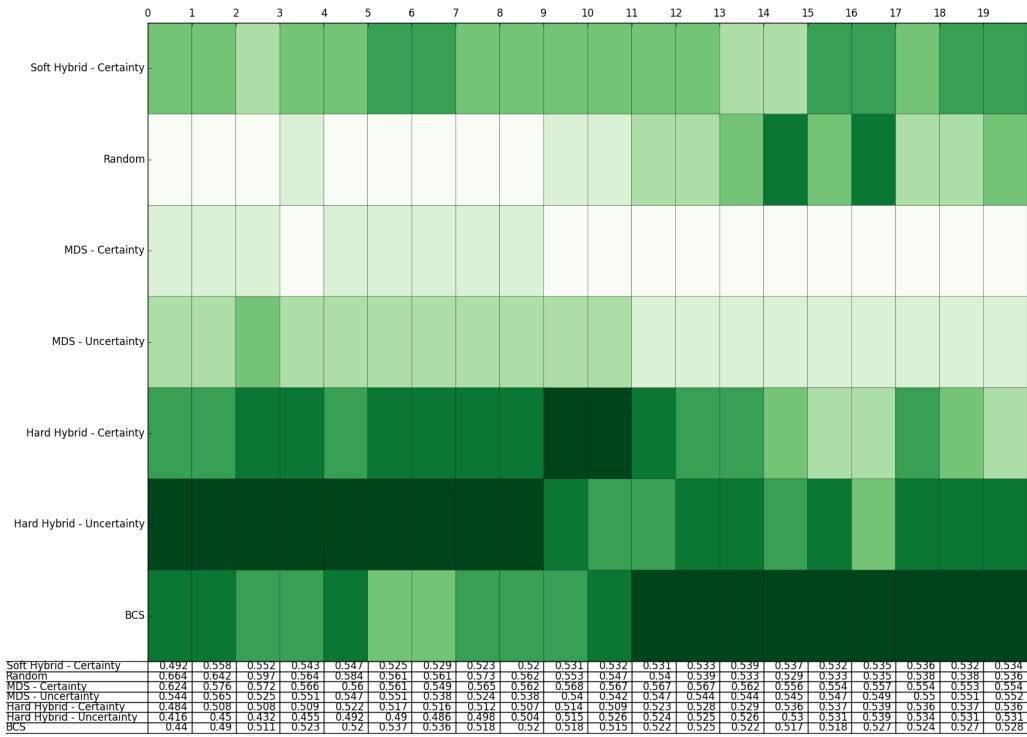


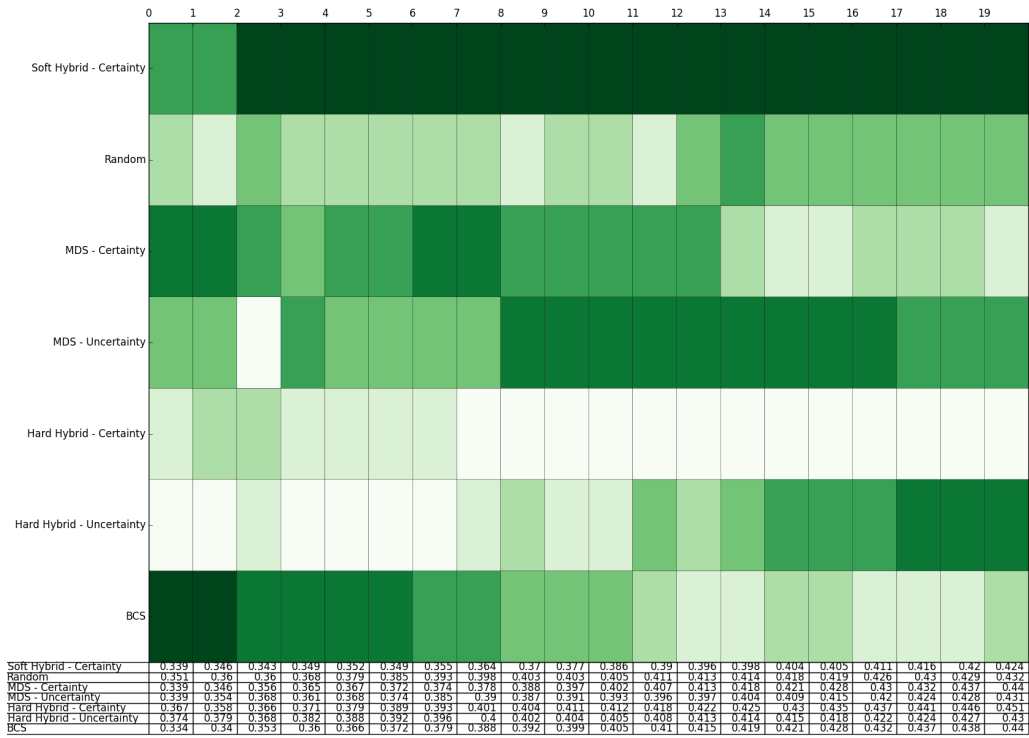
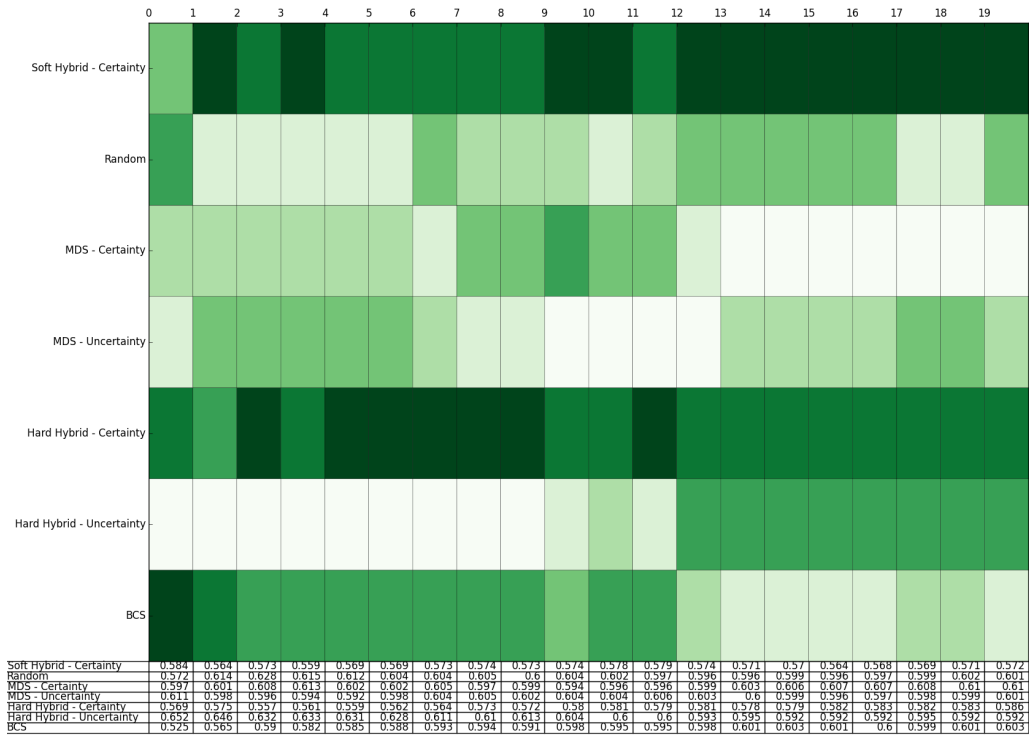




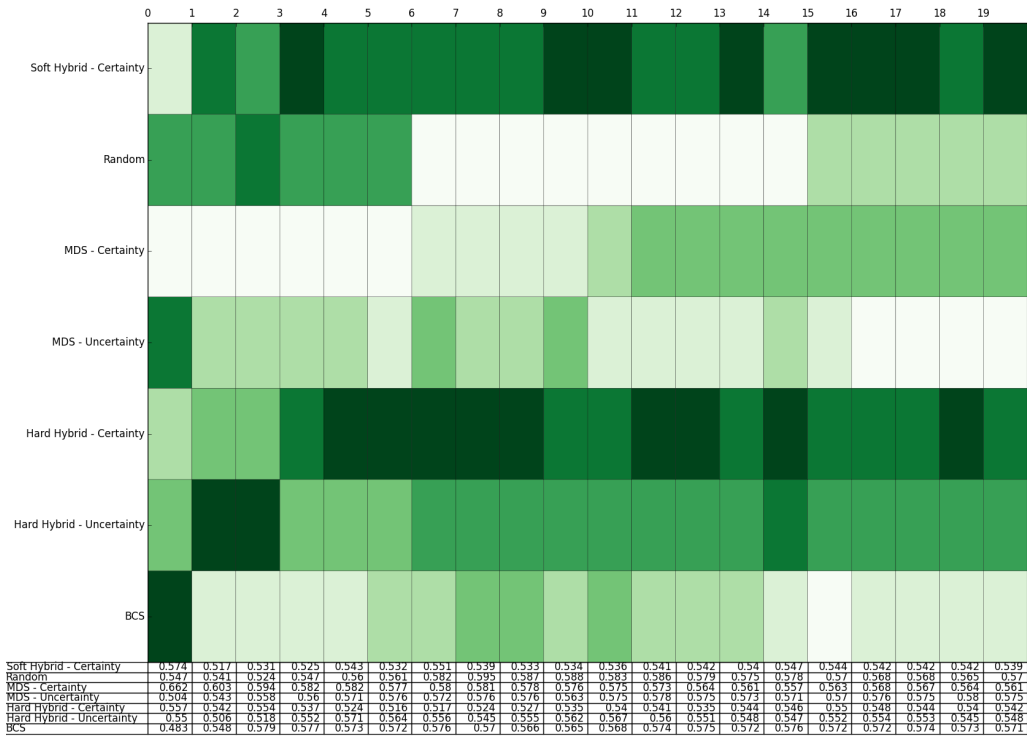
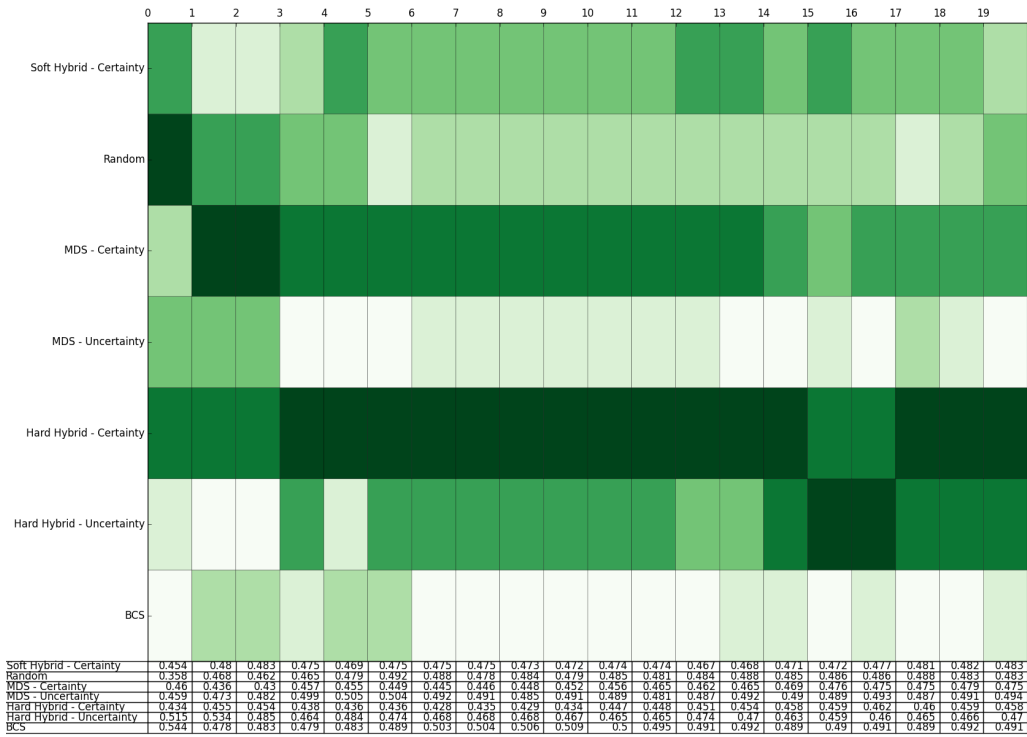


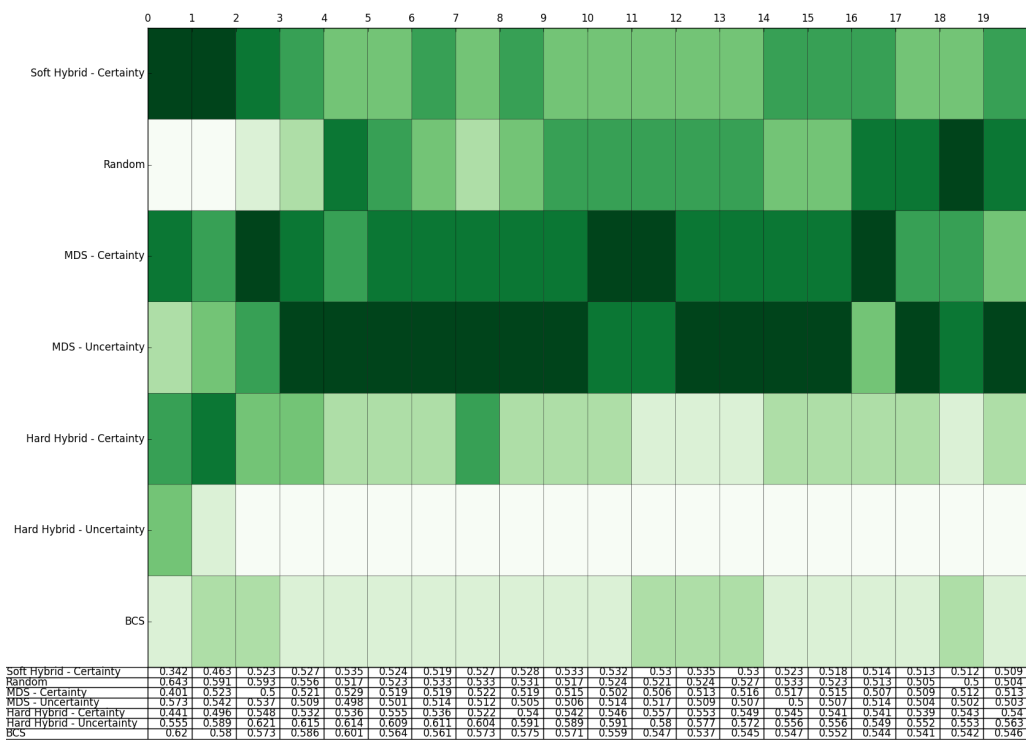
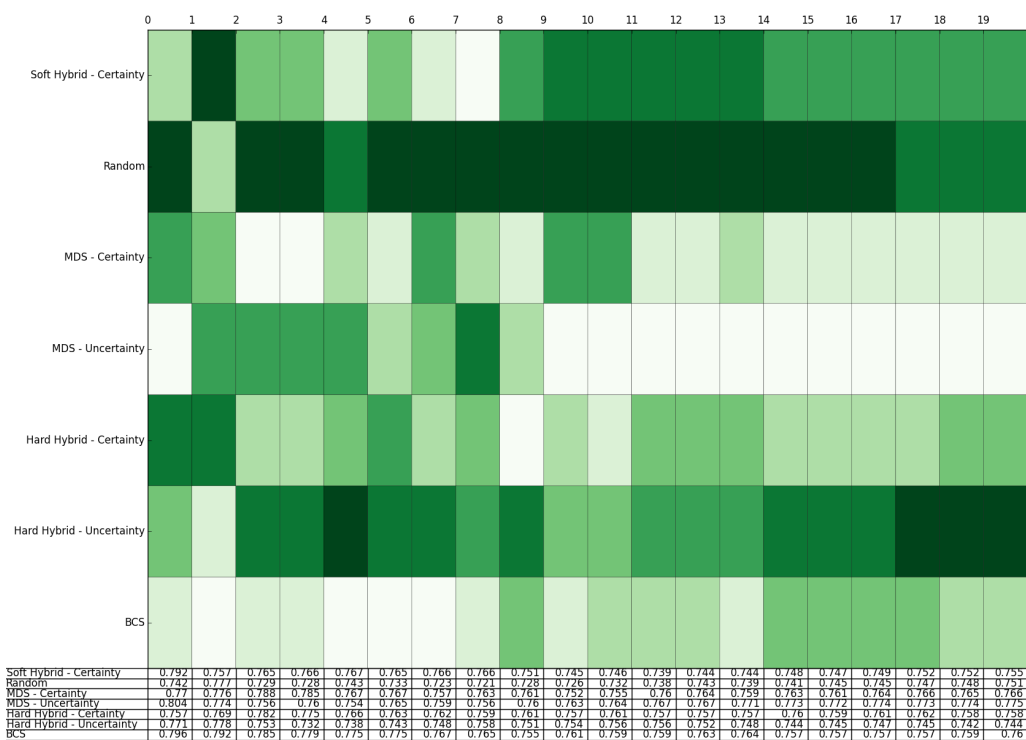


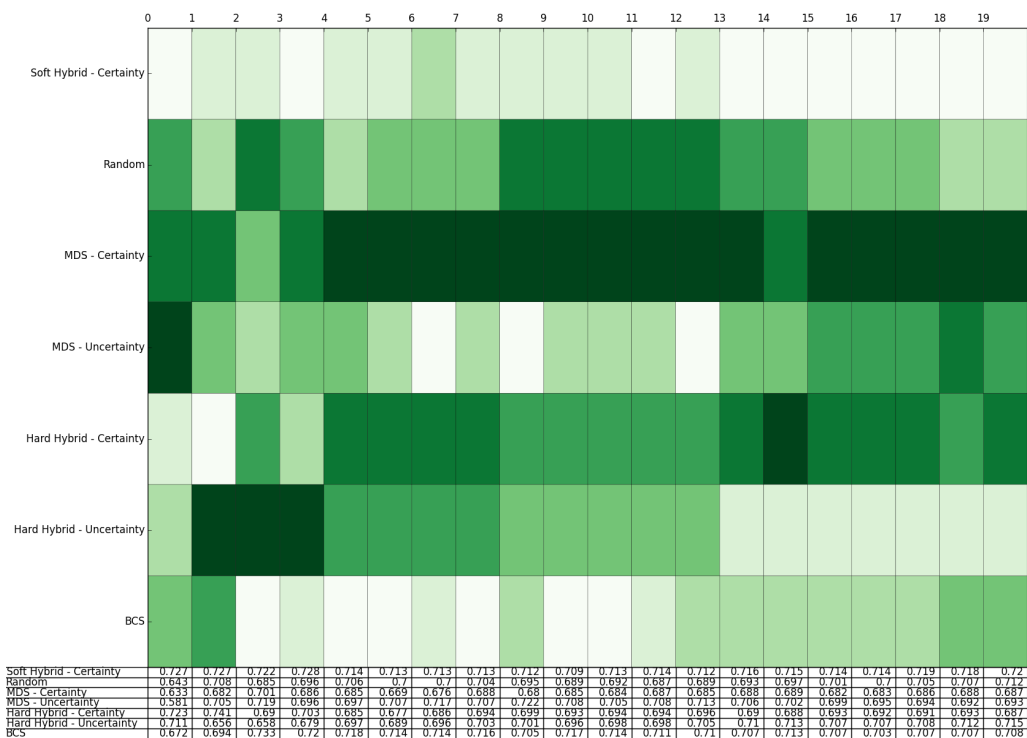
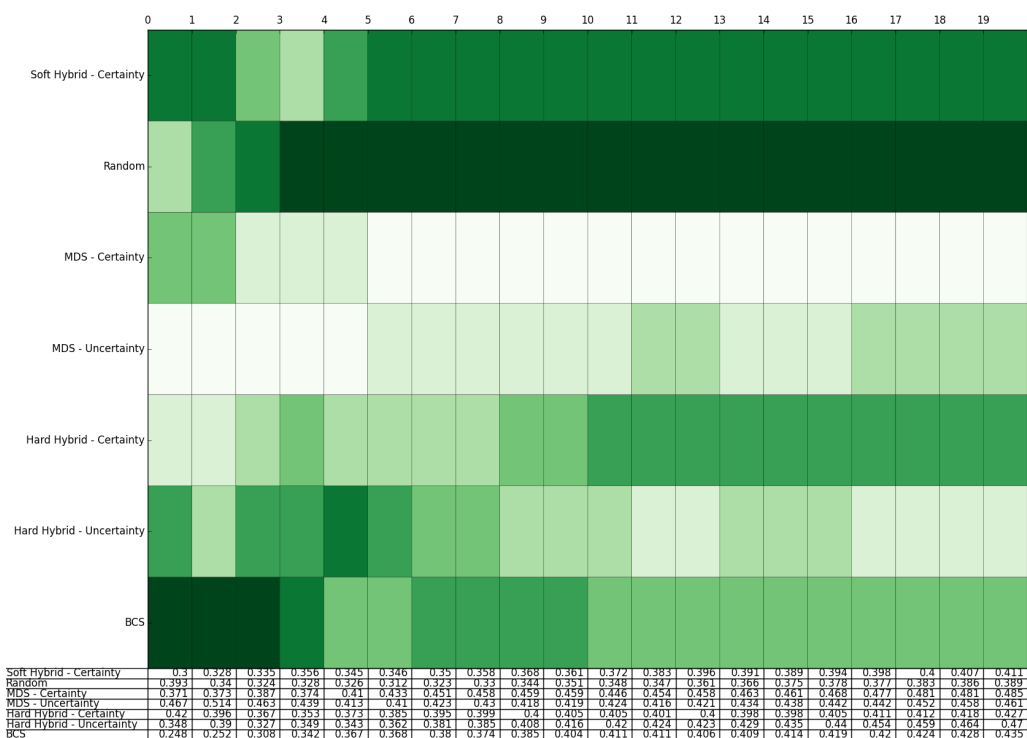


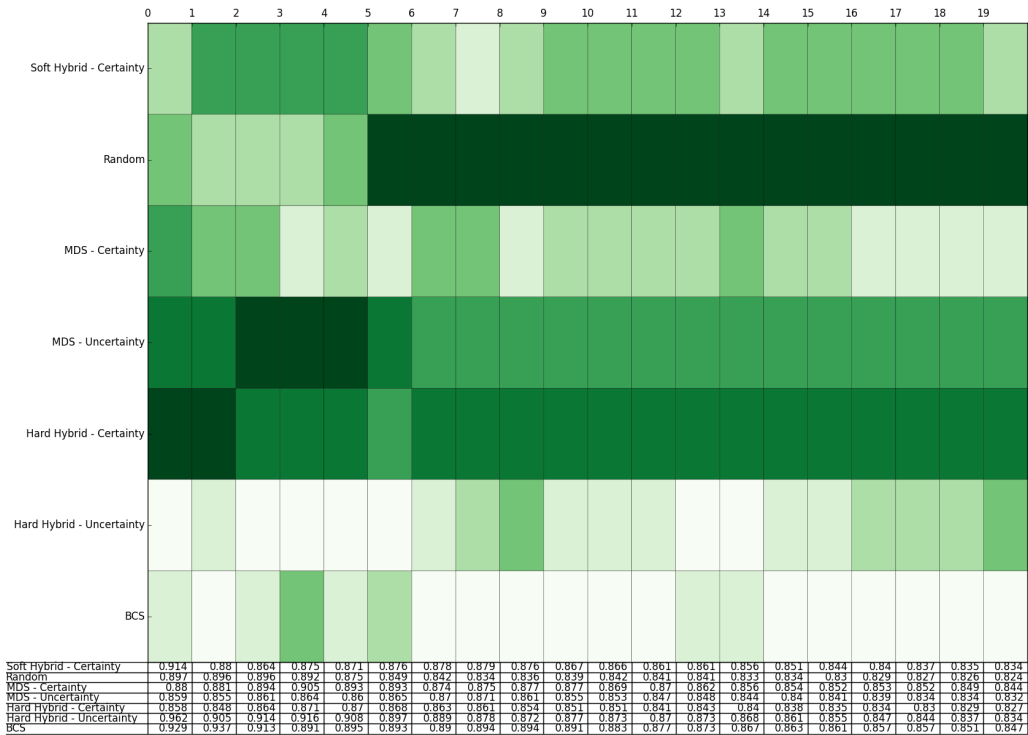
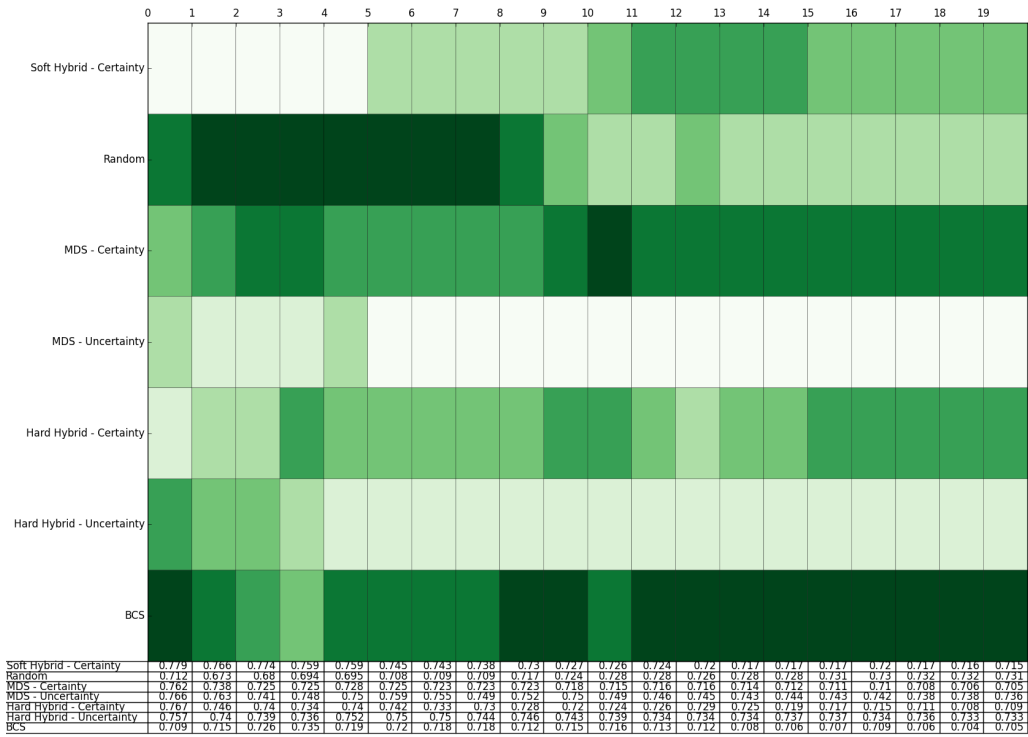


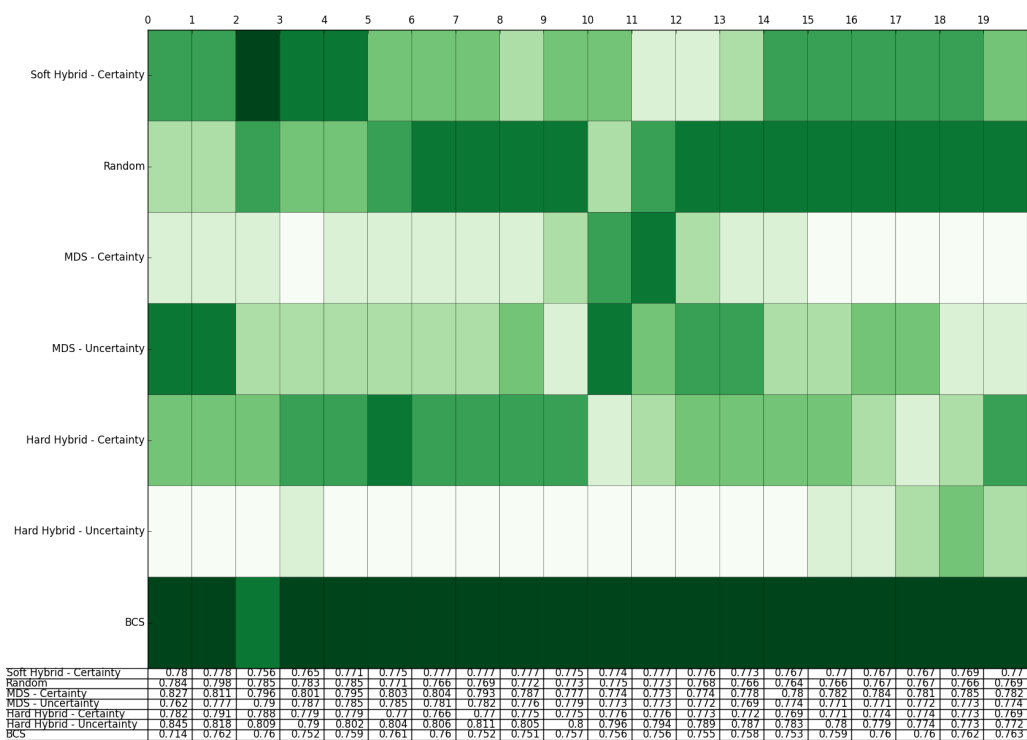




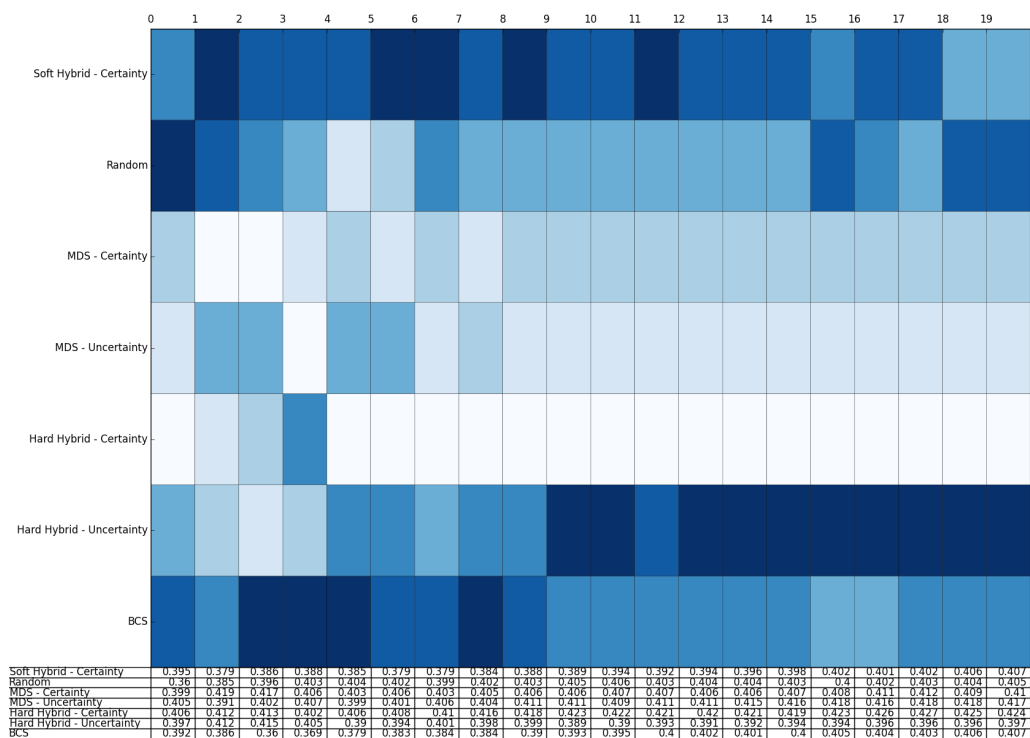


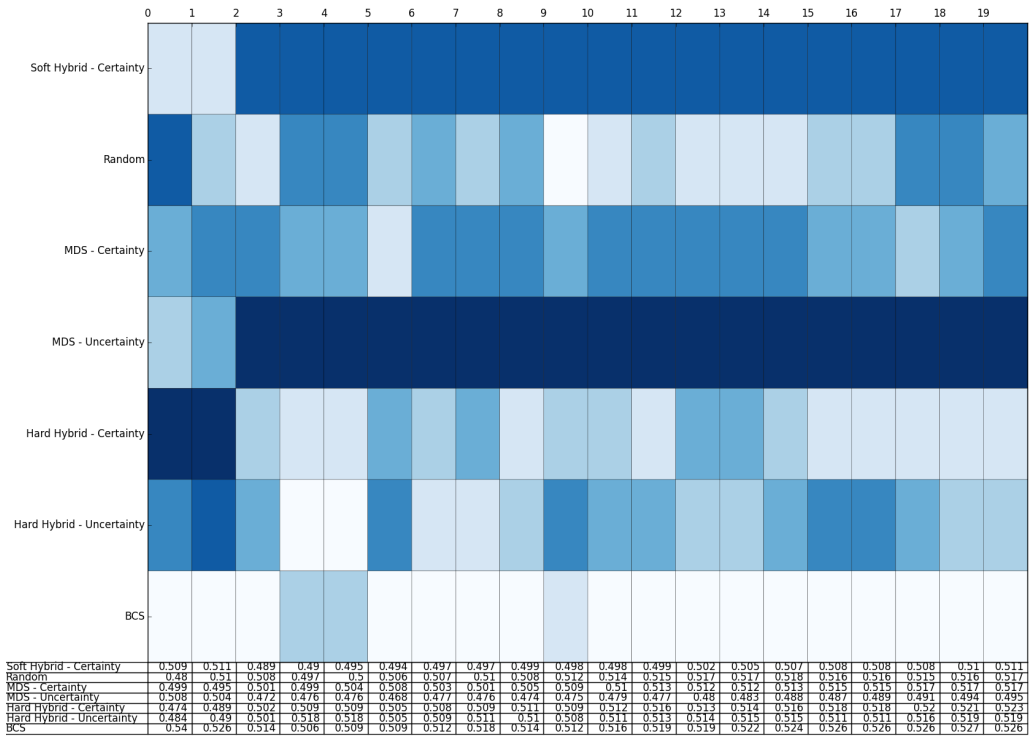
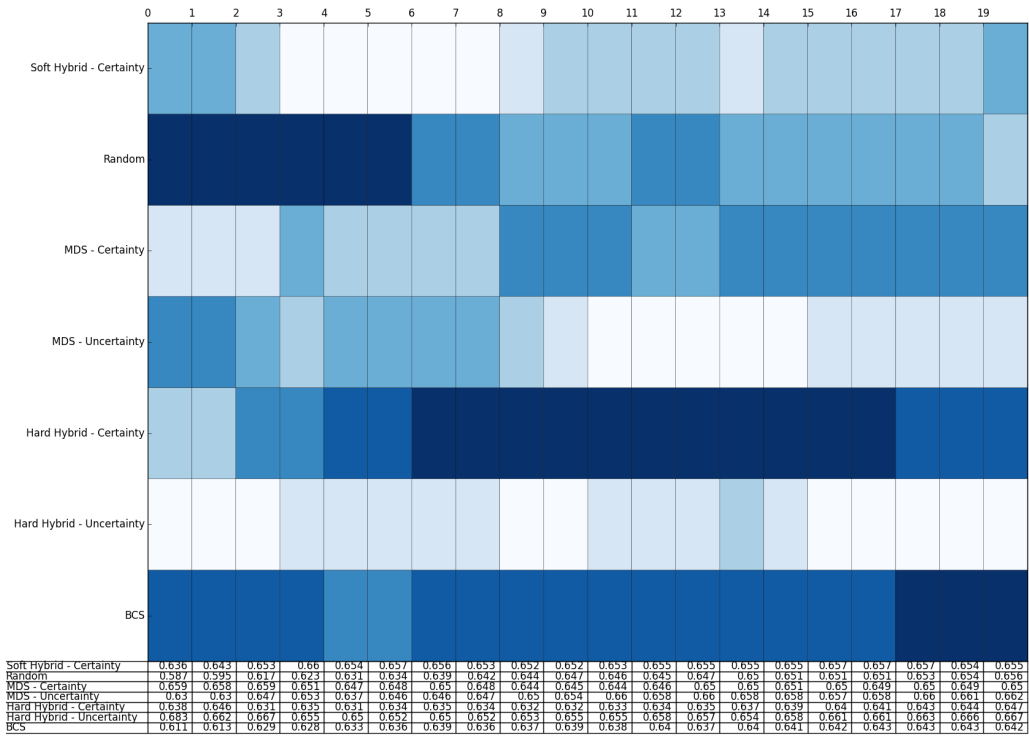


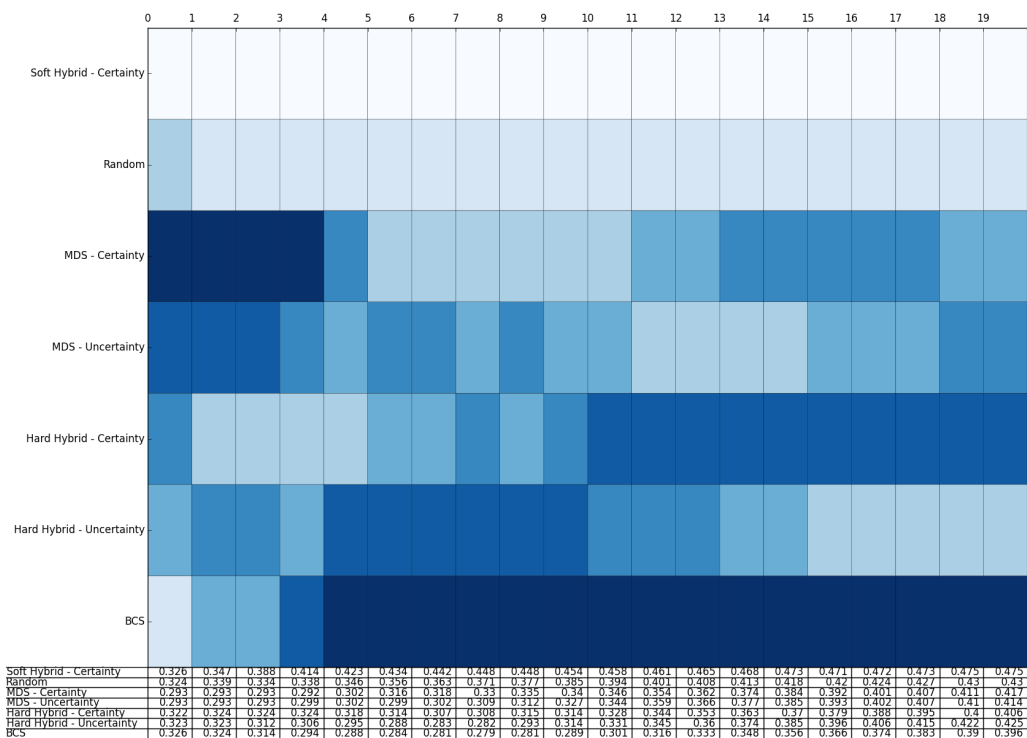
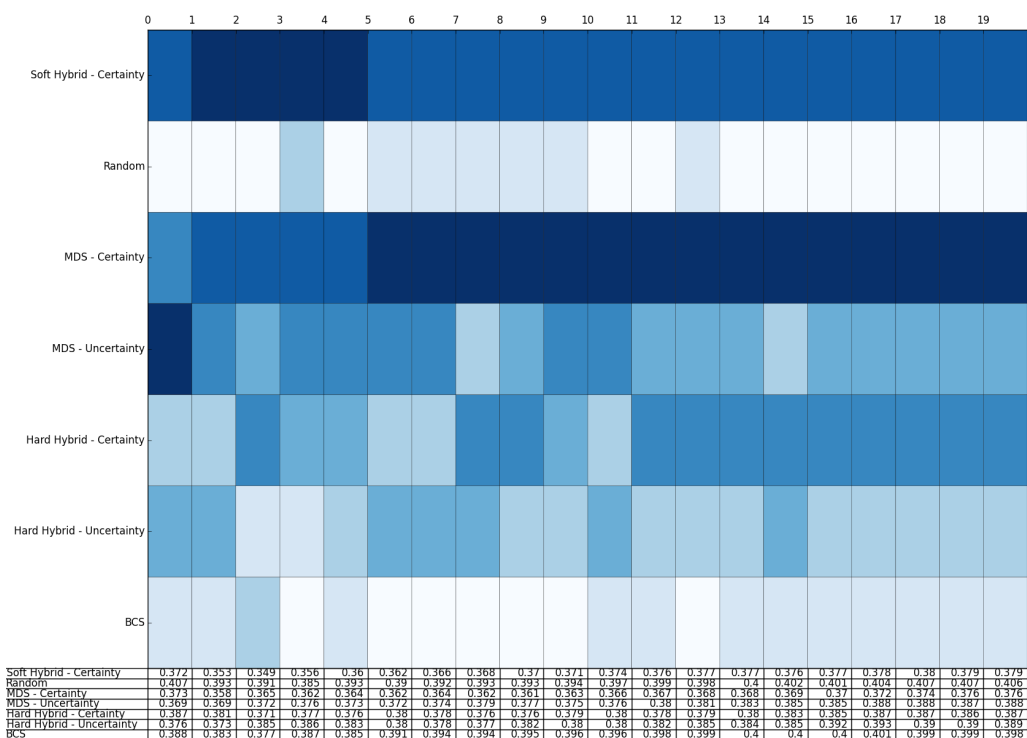


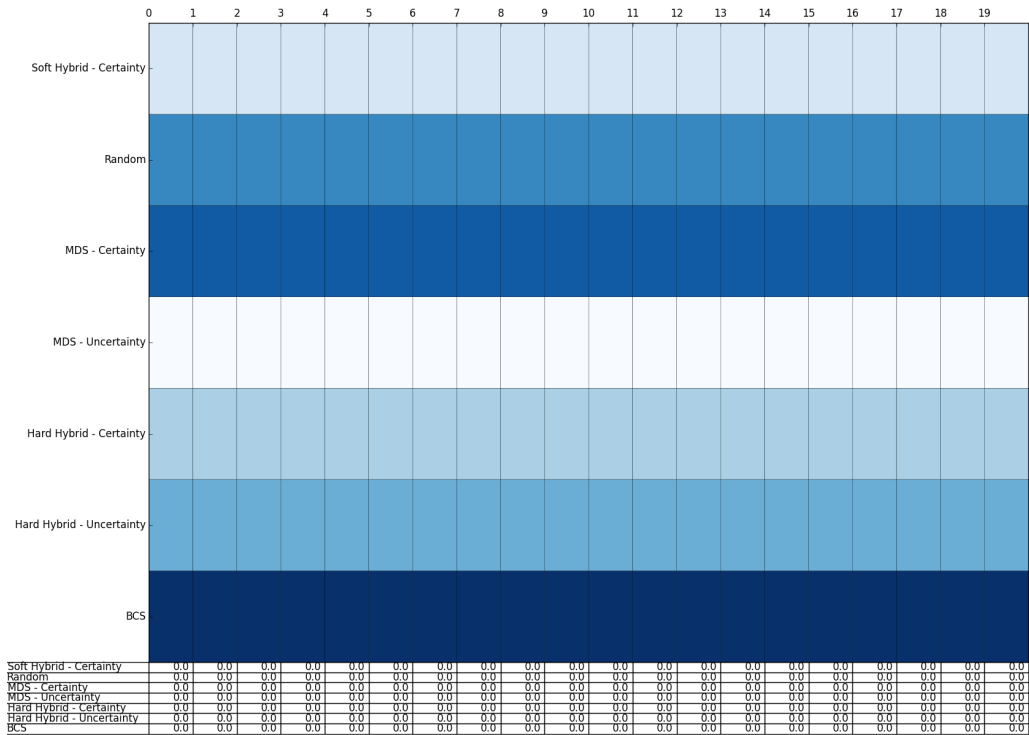
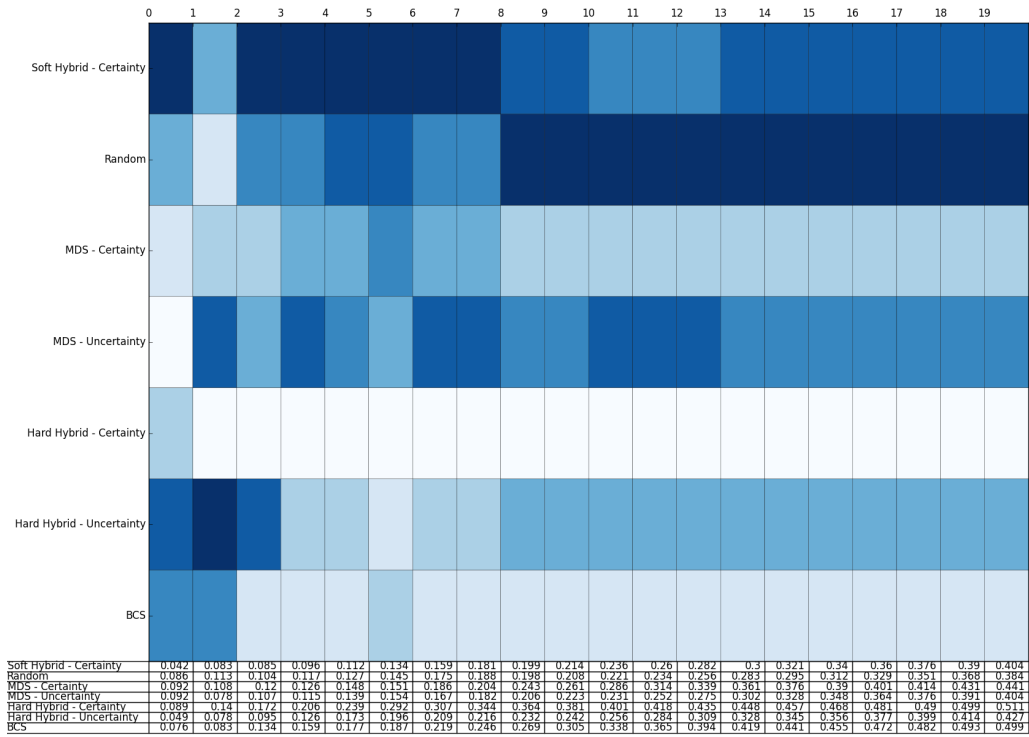


Speaker score using face annotation is given blue heat maps.

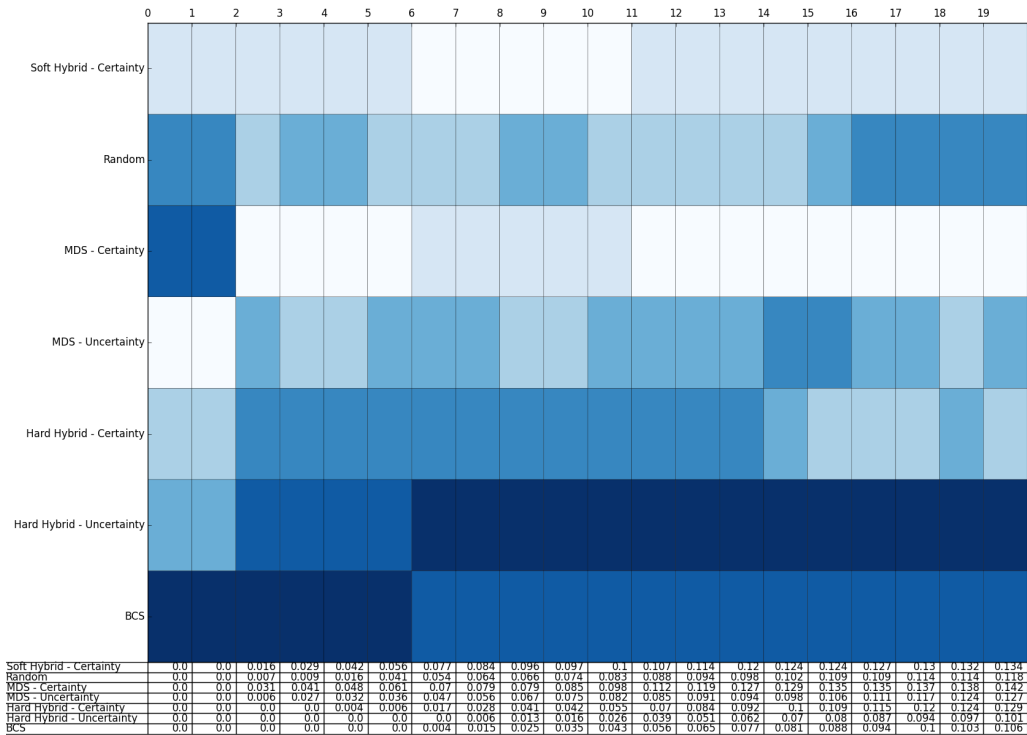
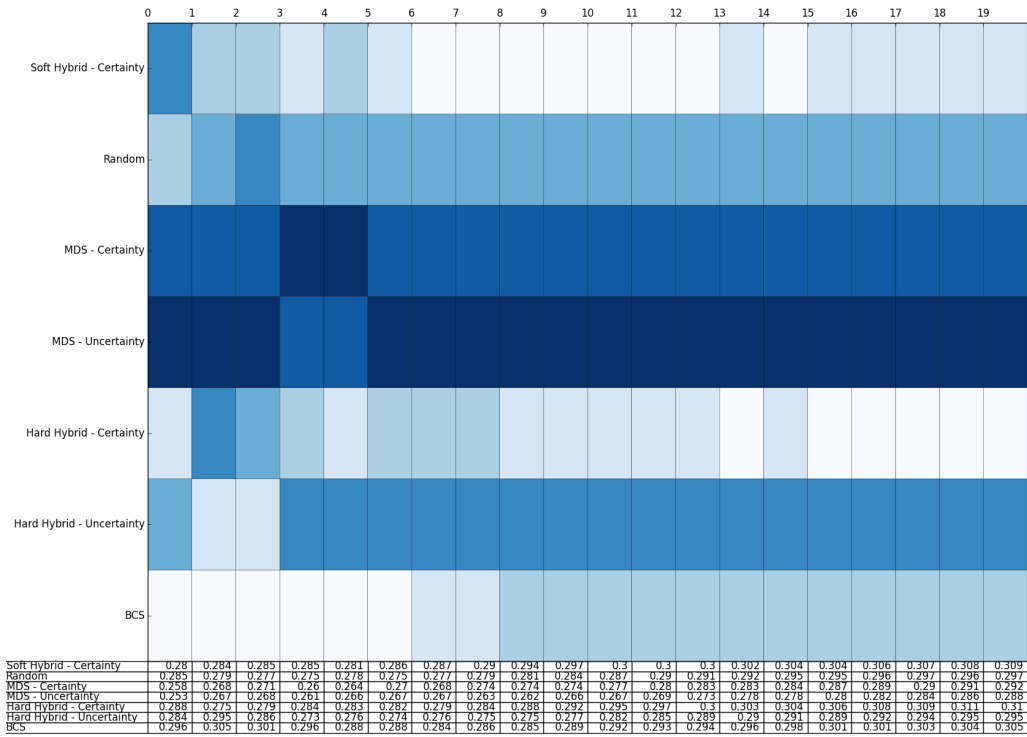


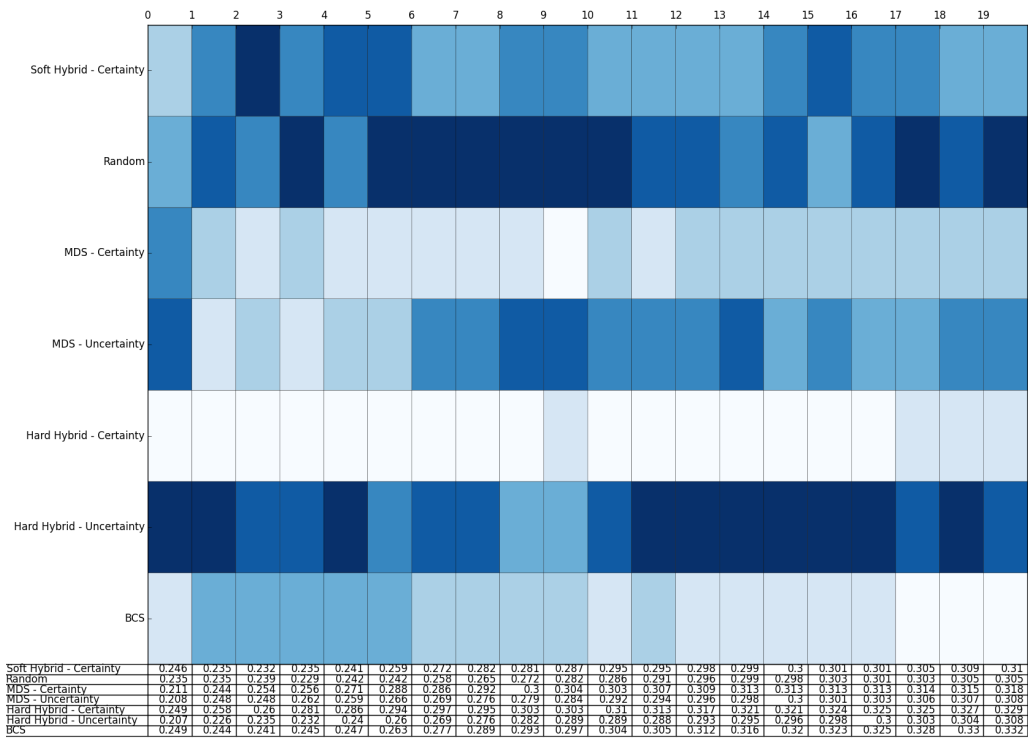
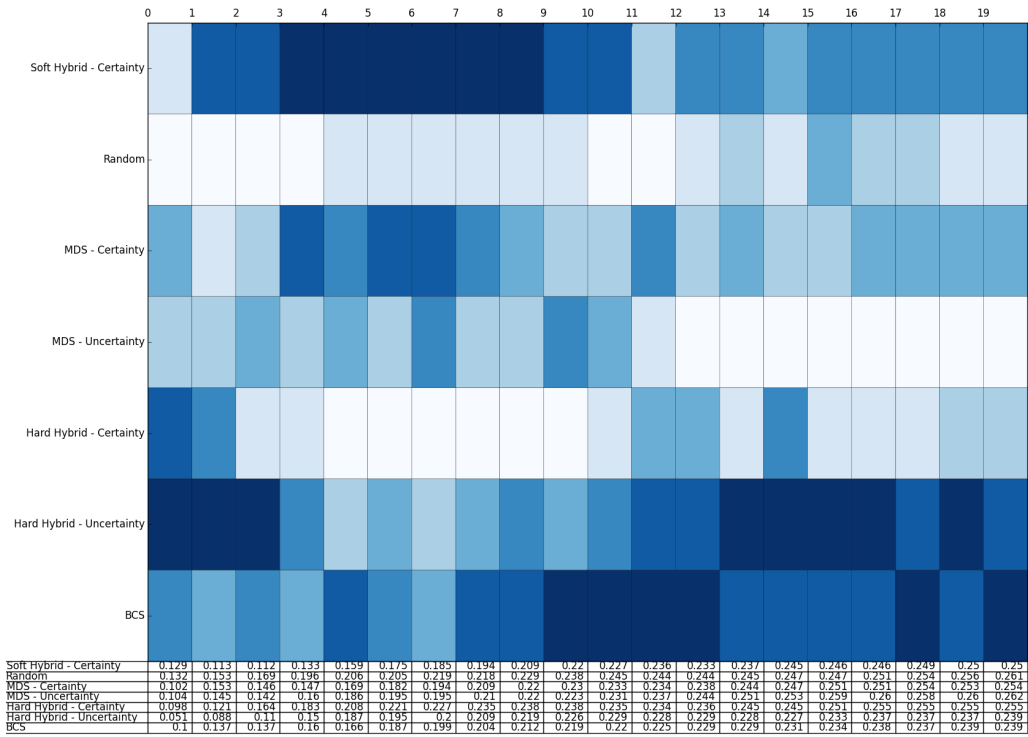


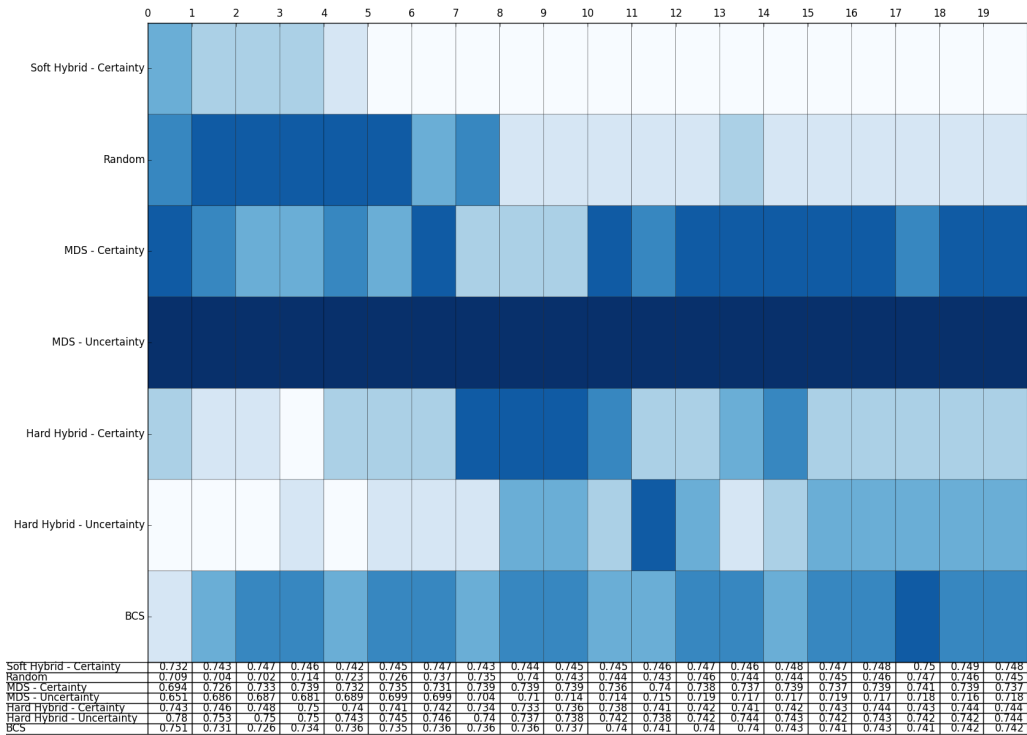
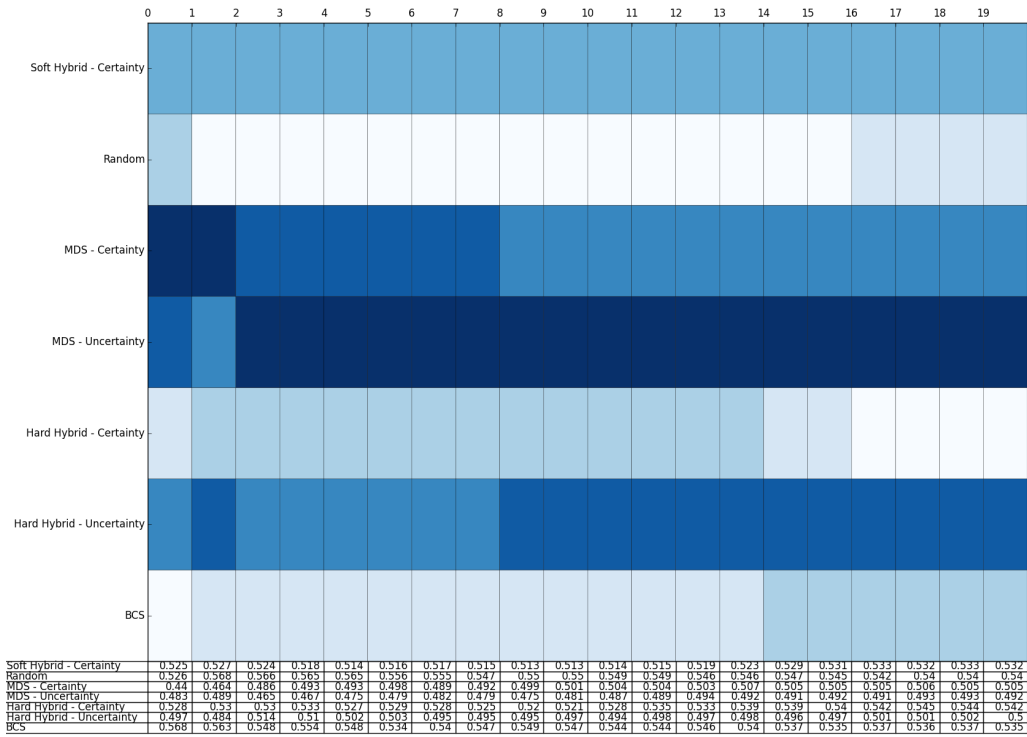


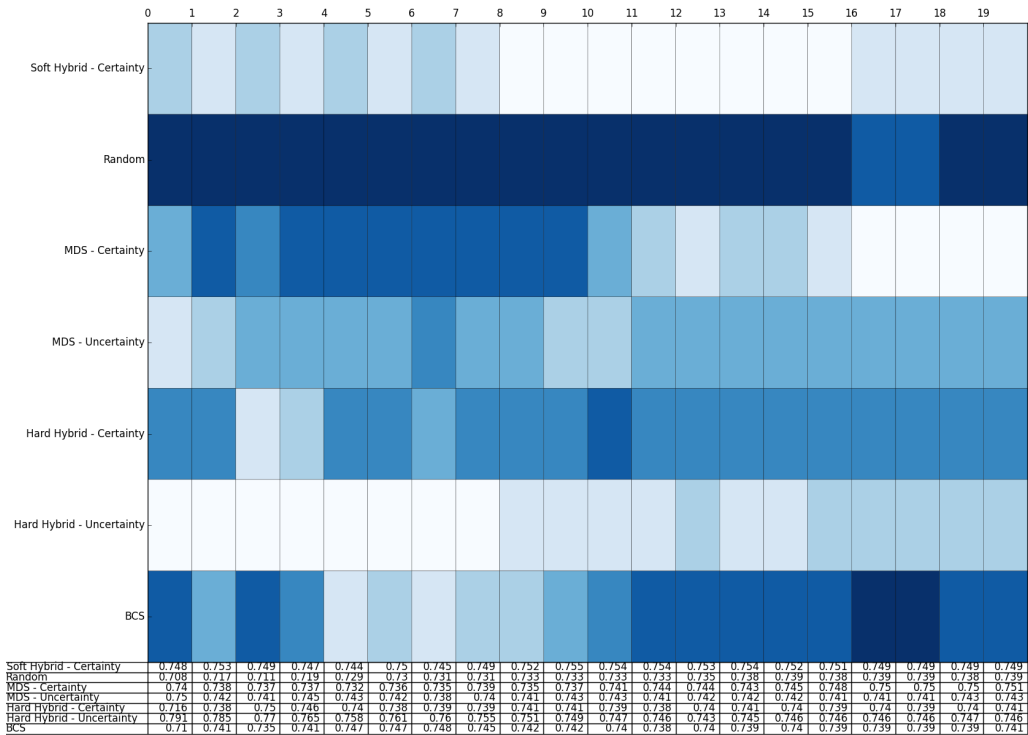
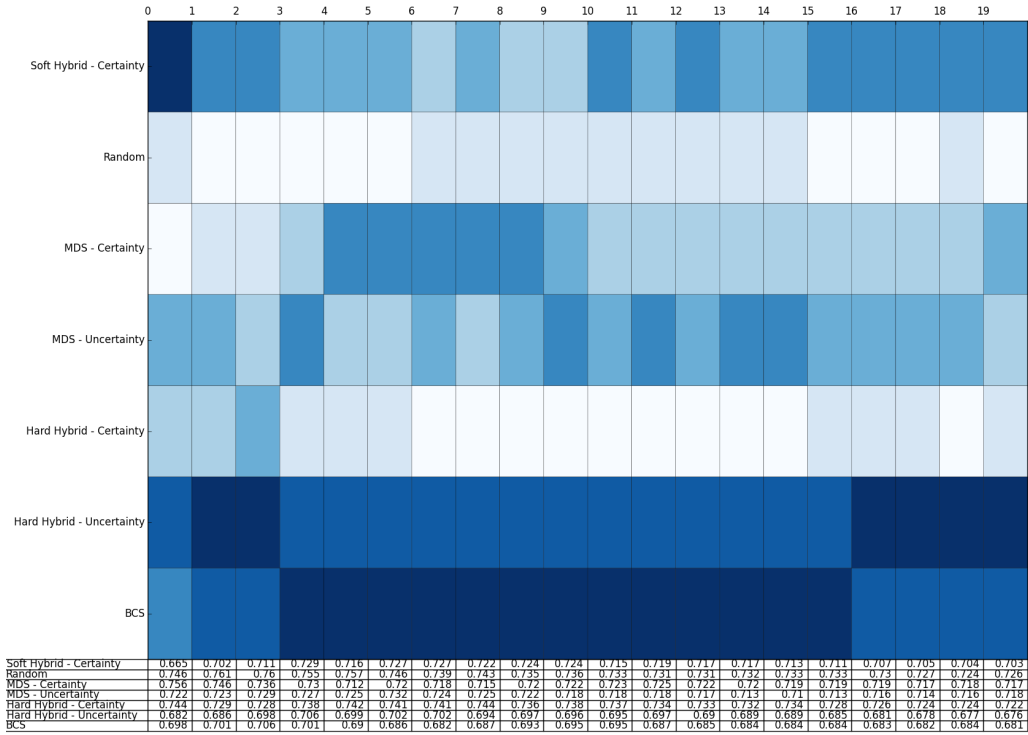


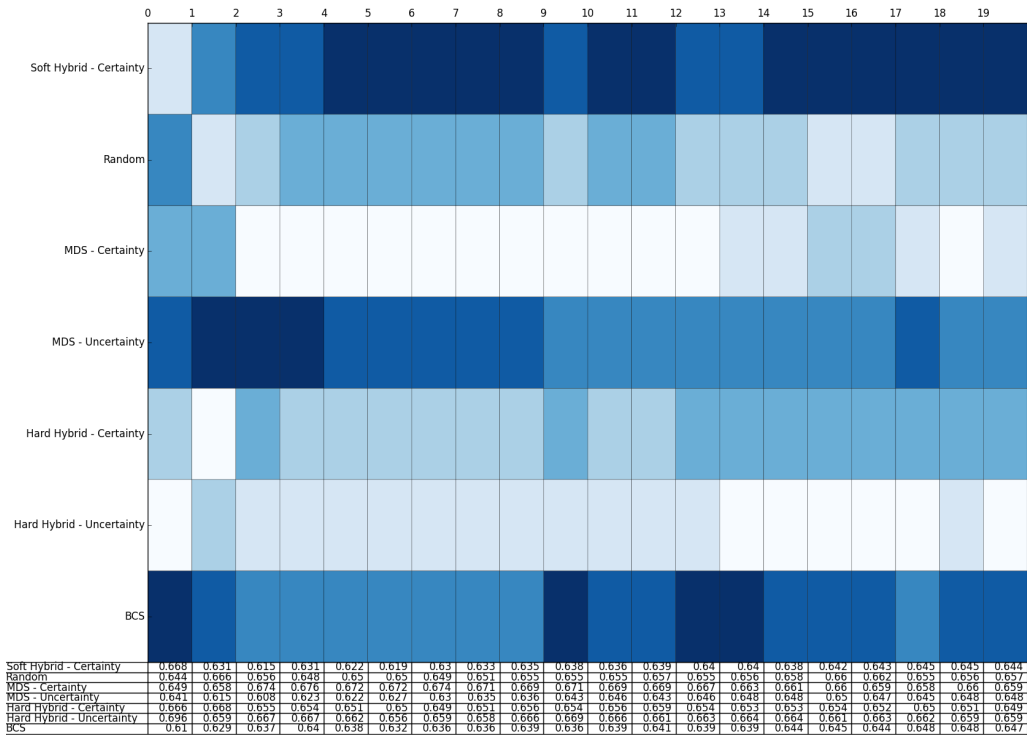
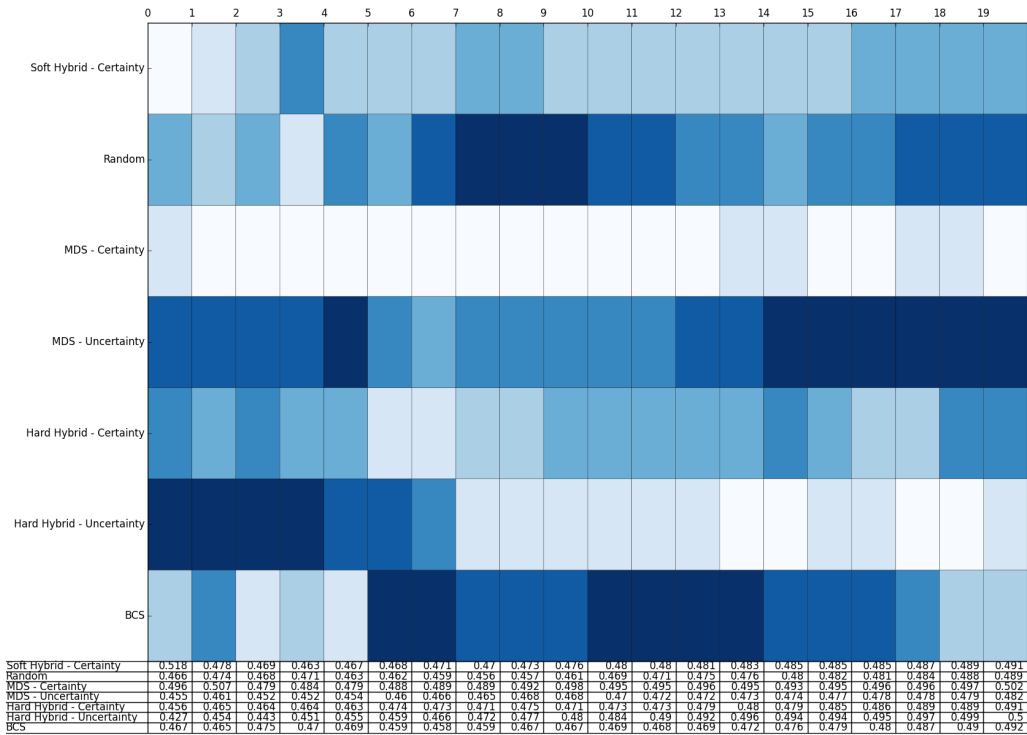


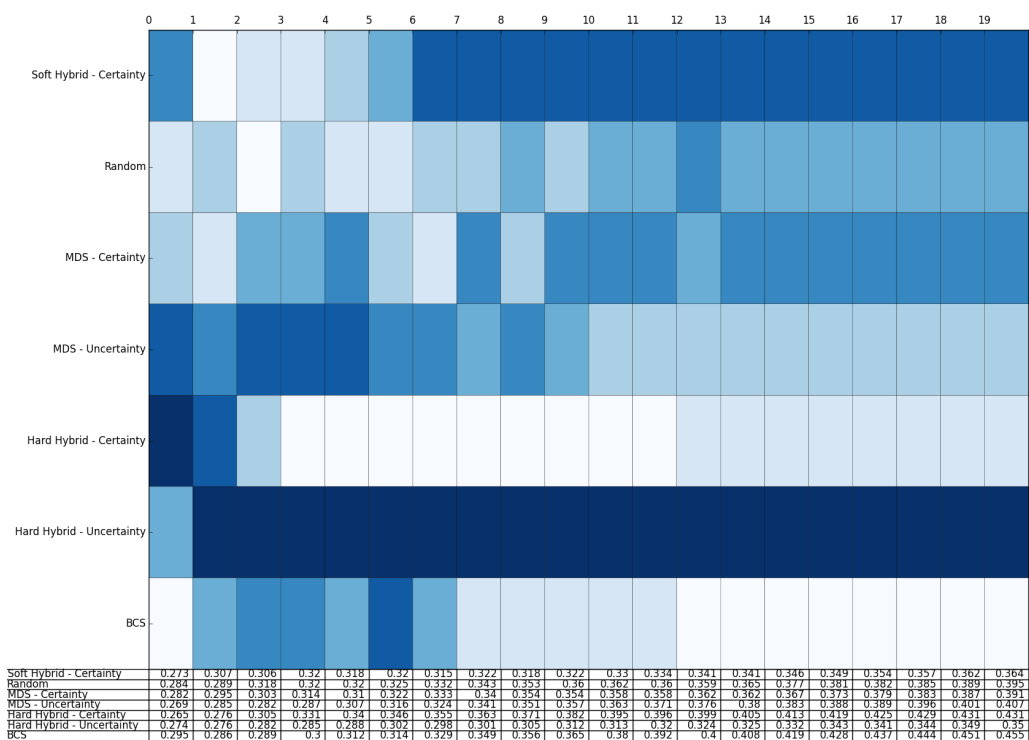
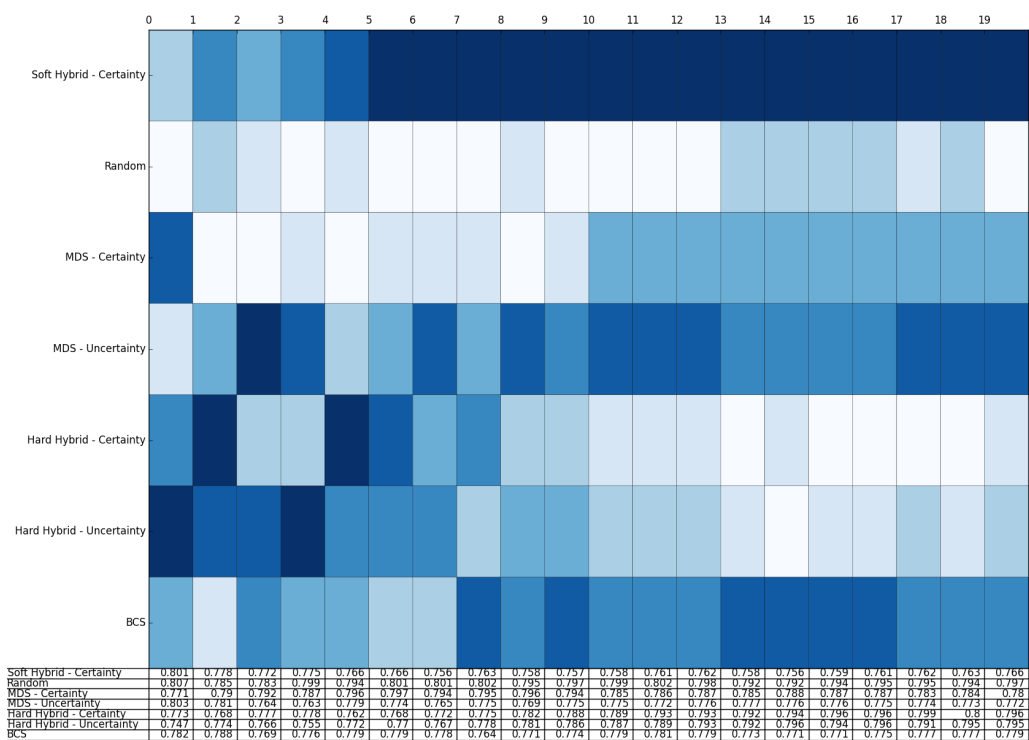


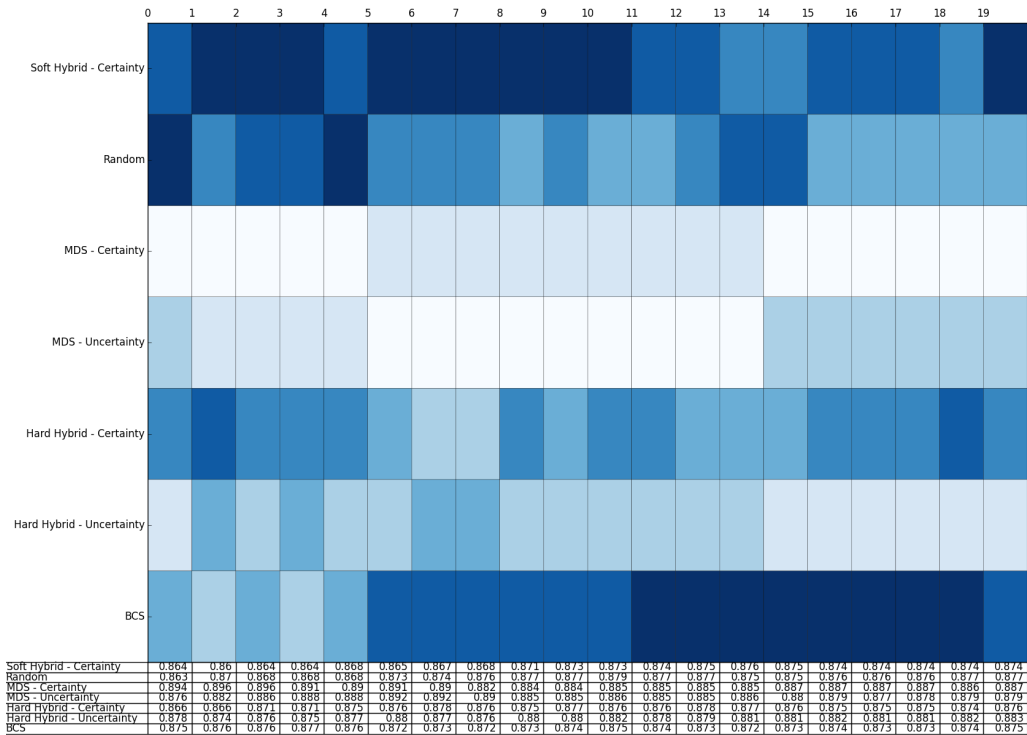
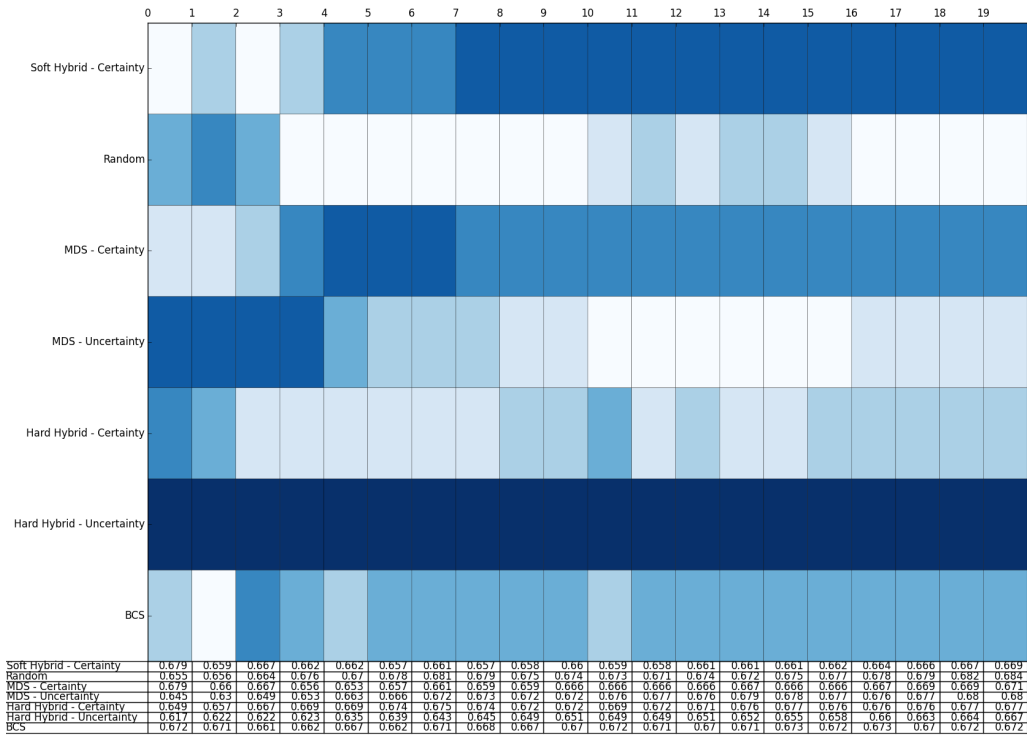


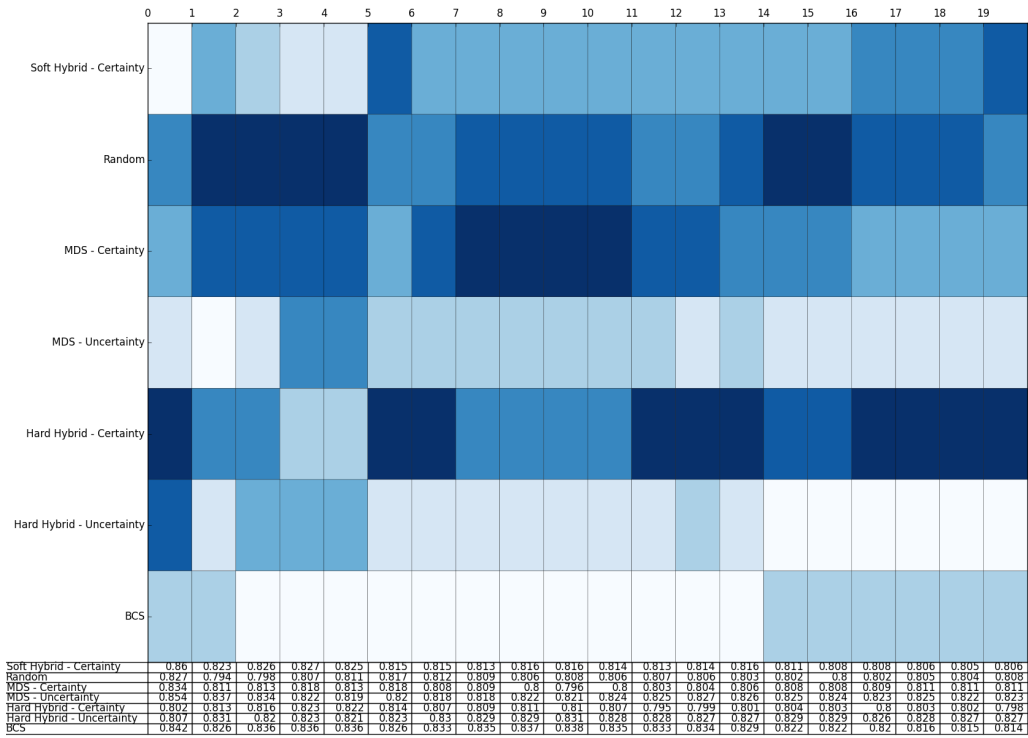
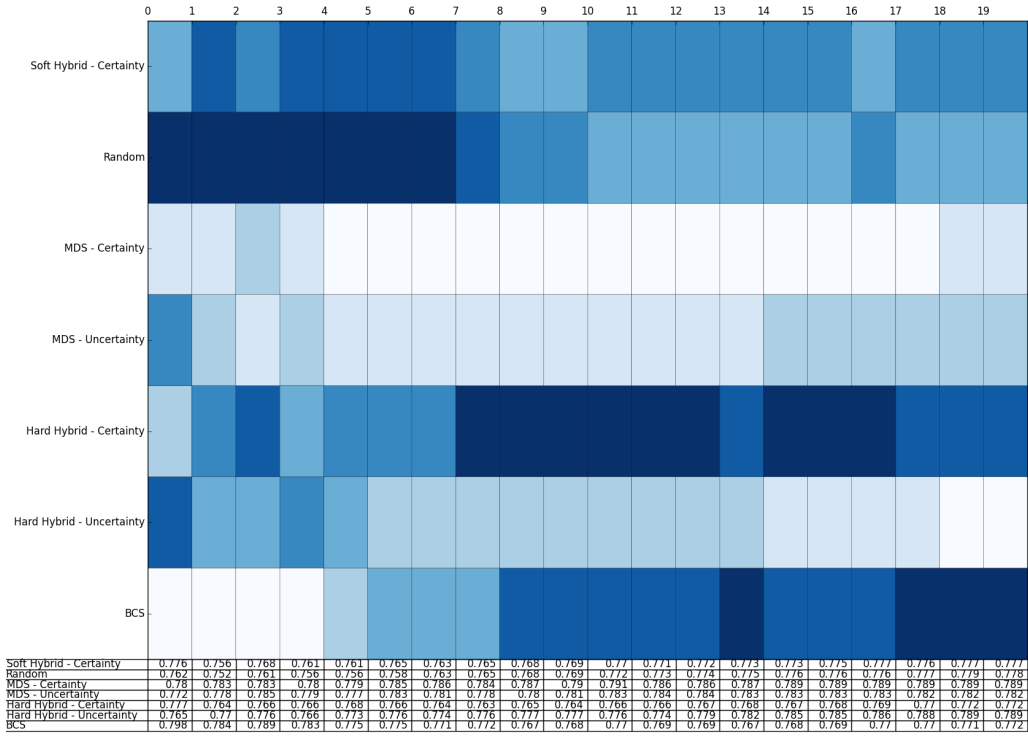




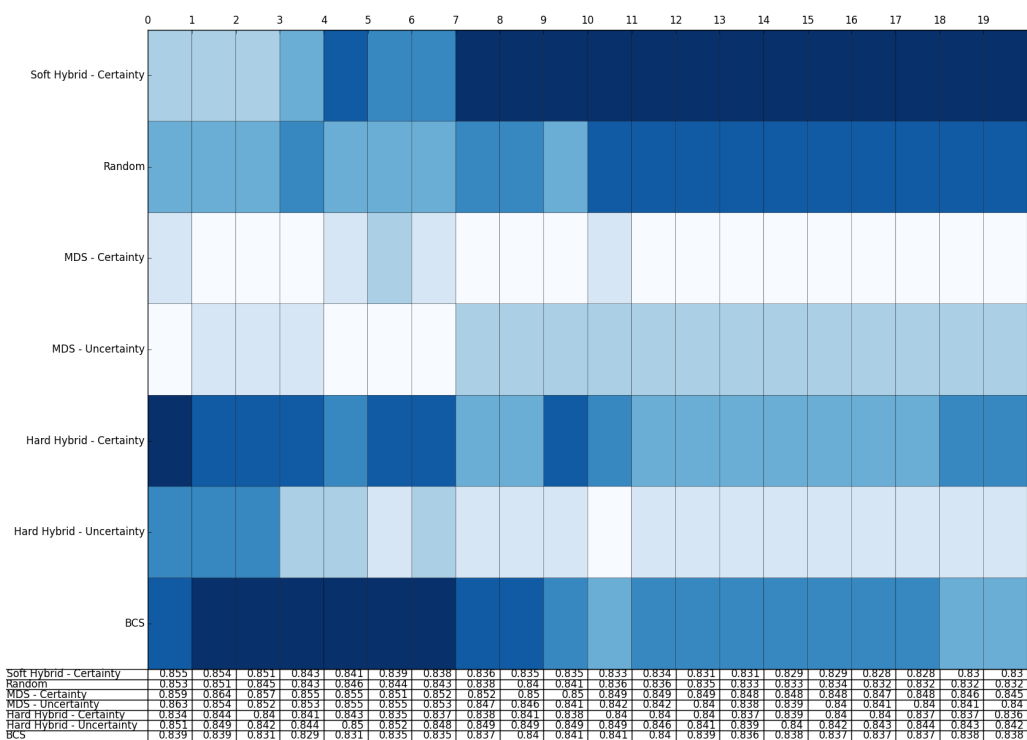
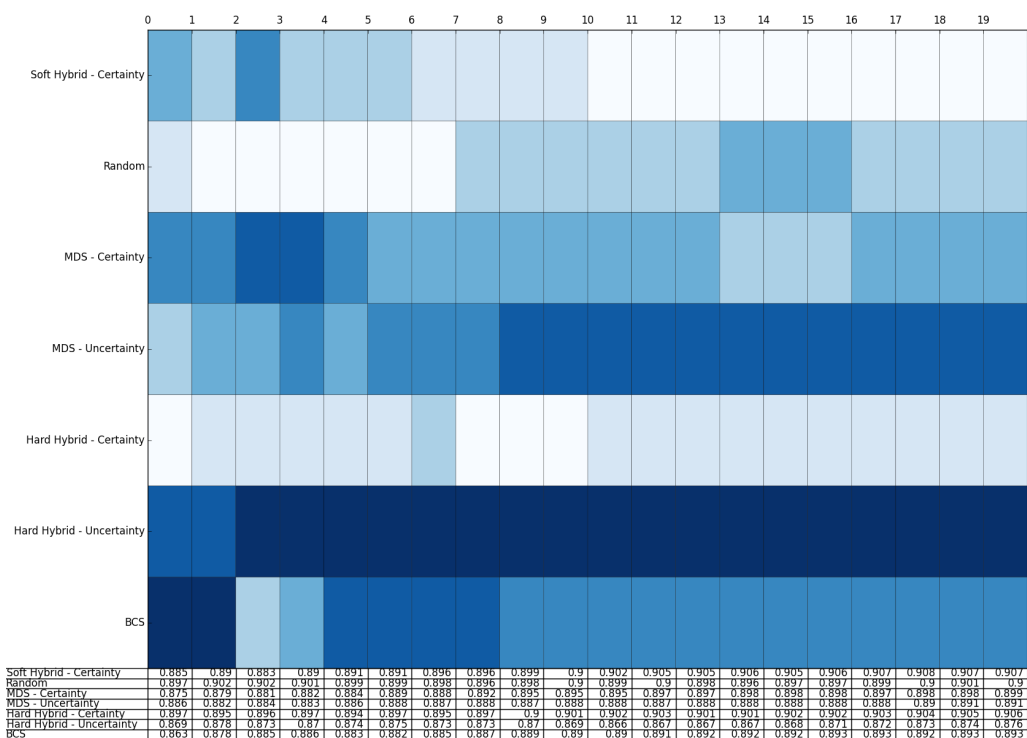


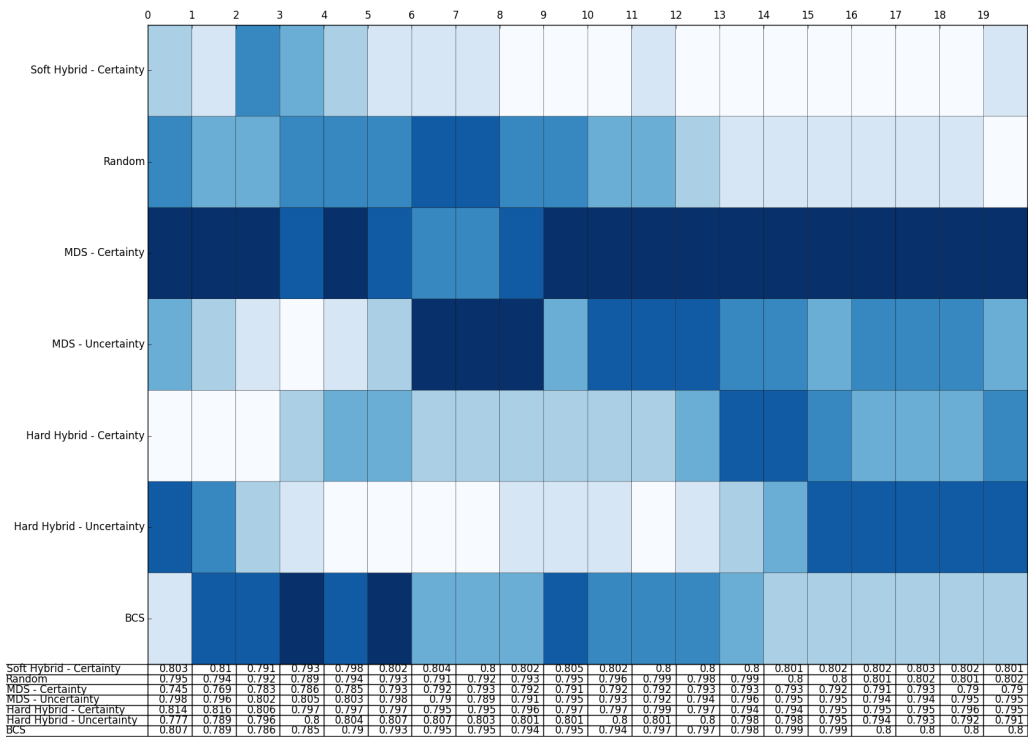
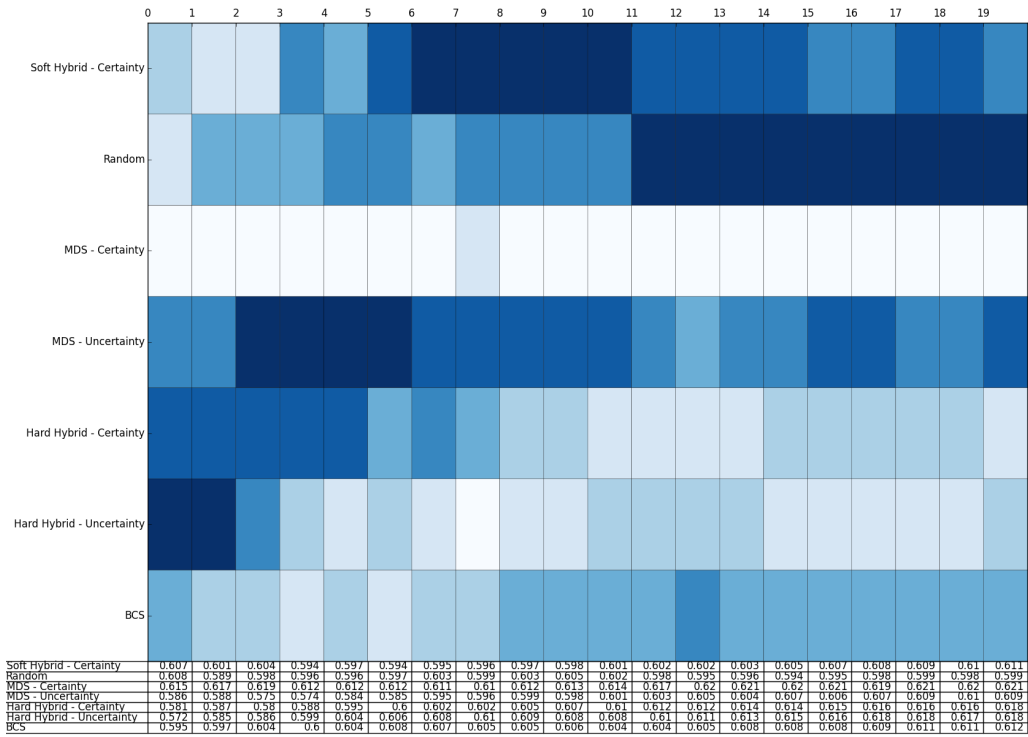


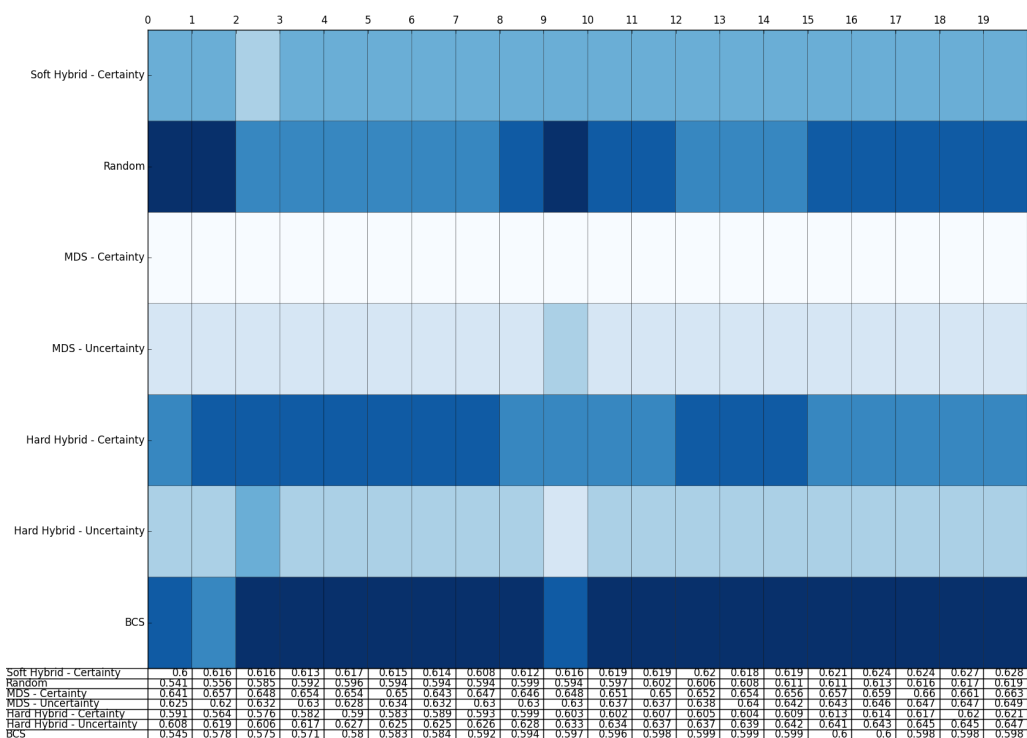




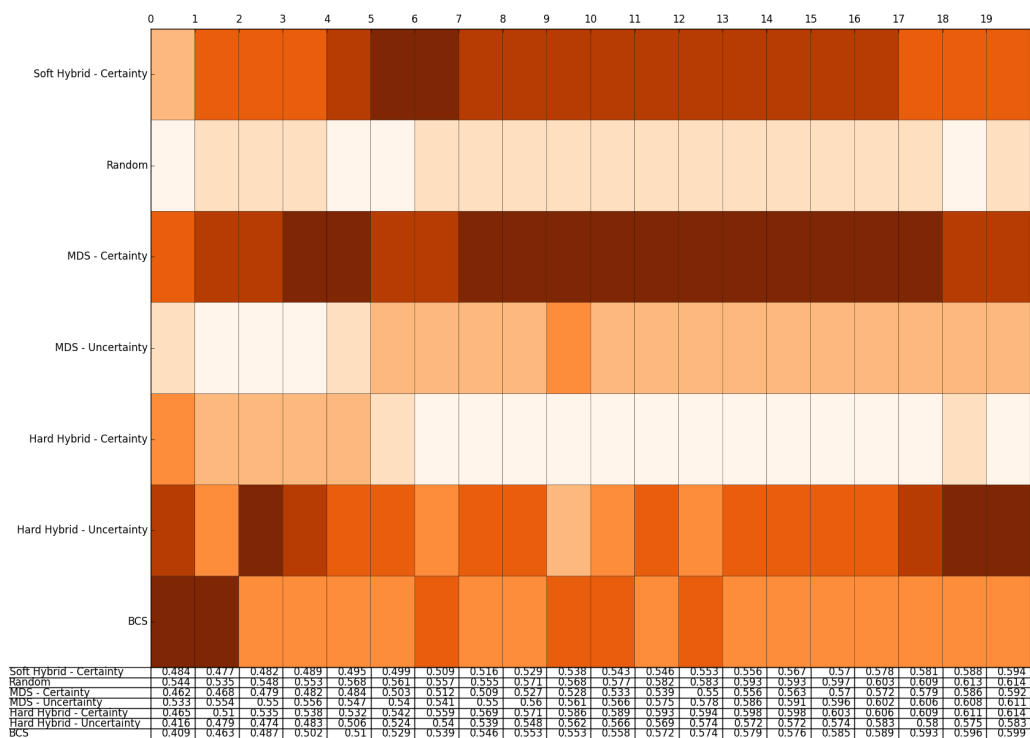


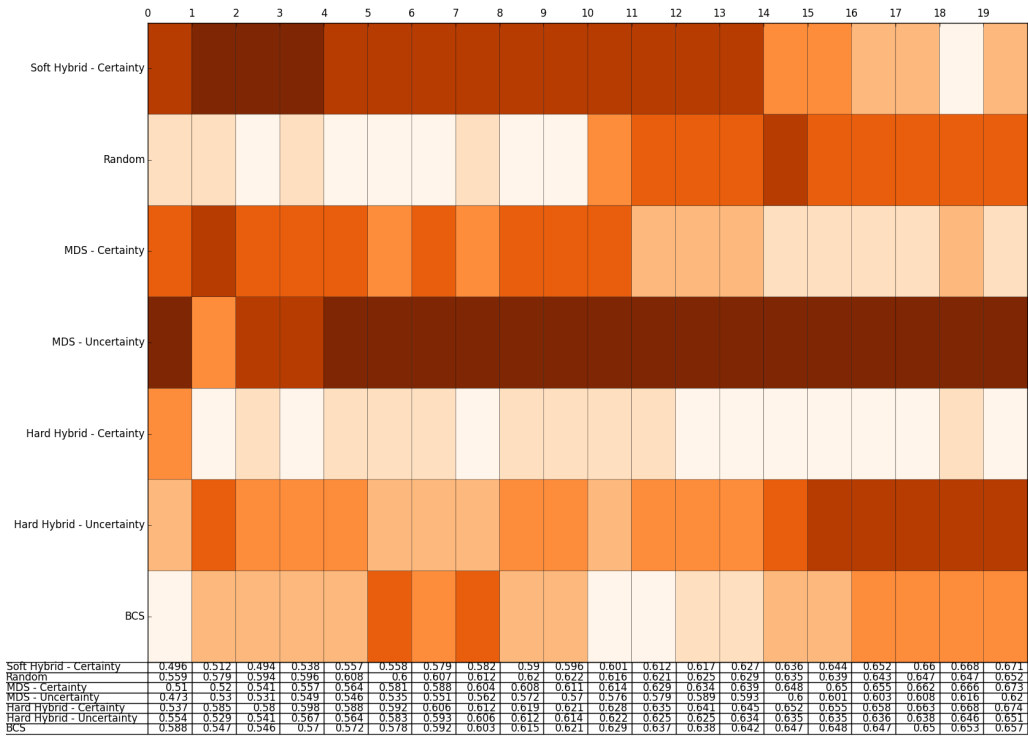
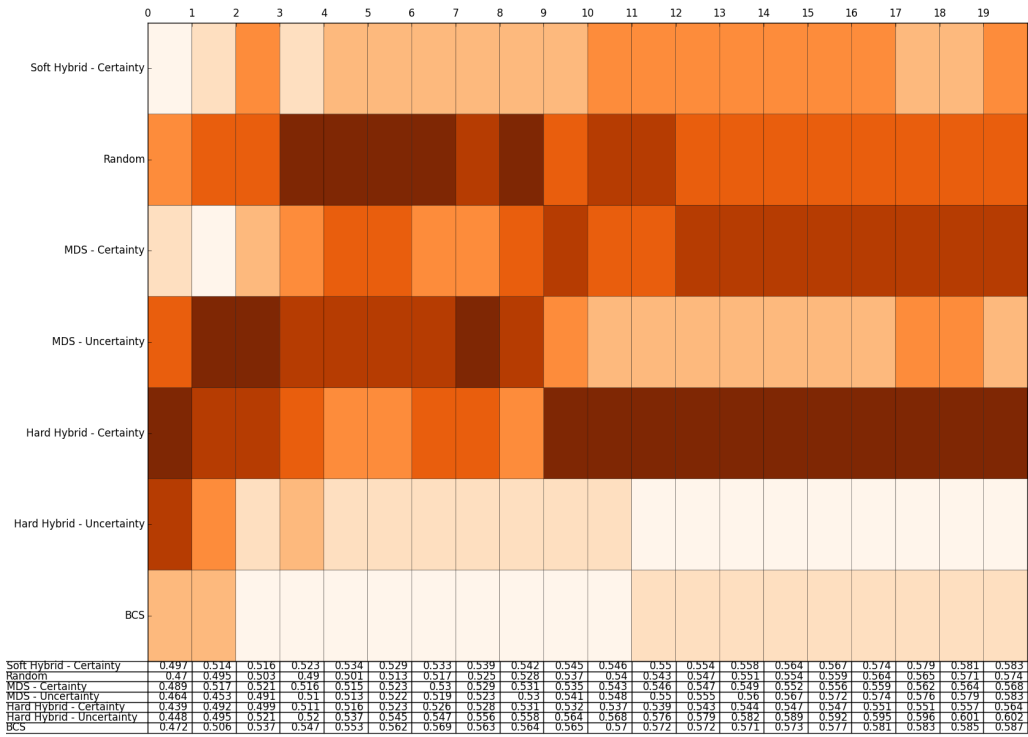


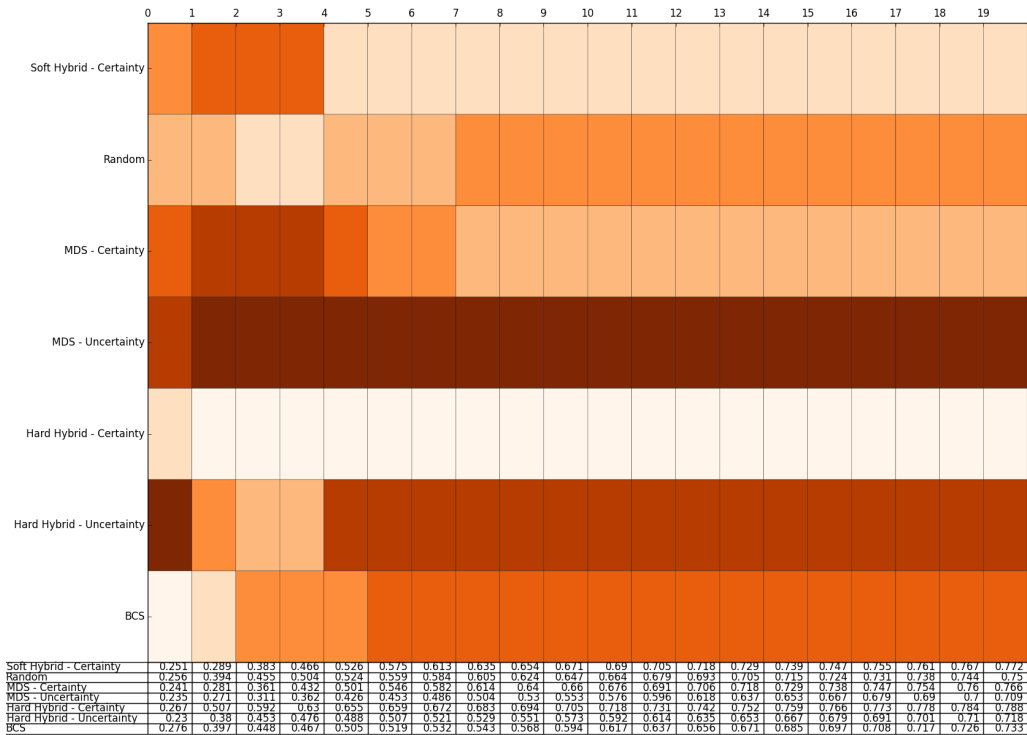
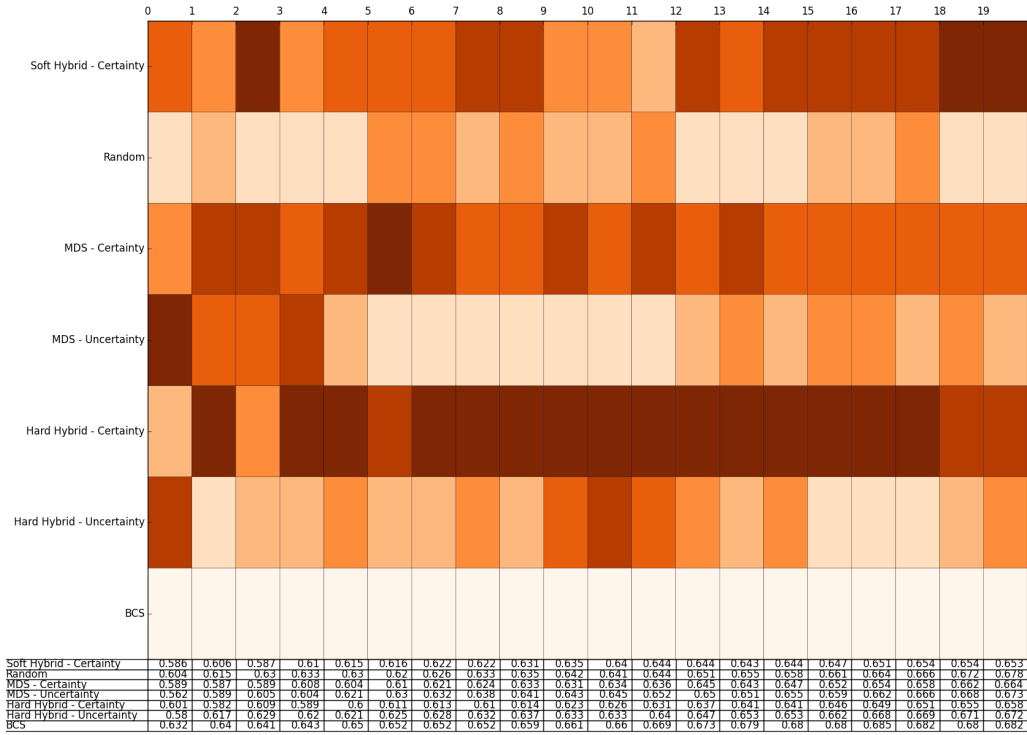


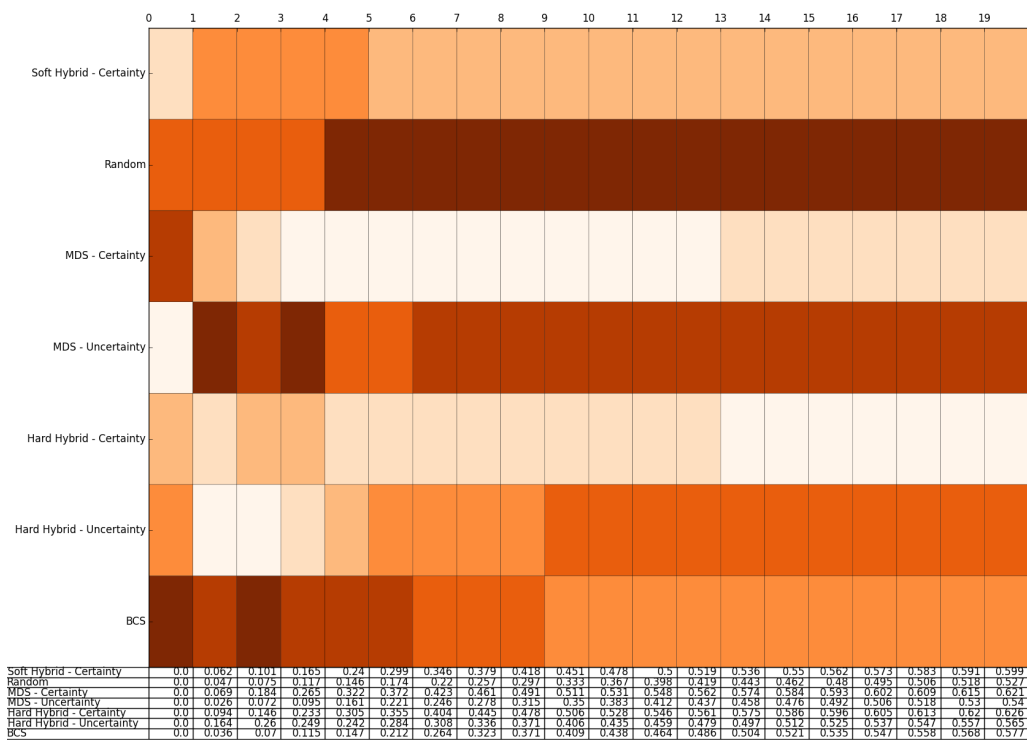
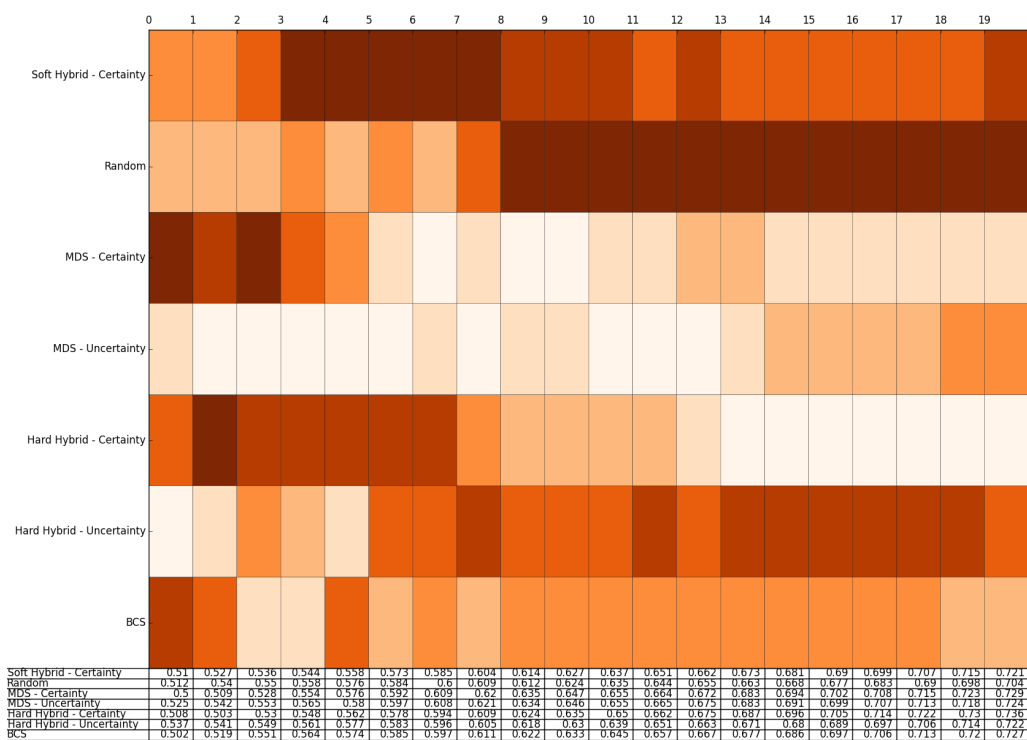


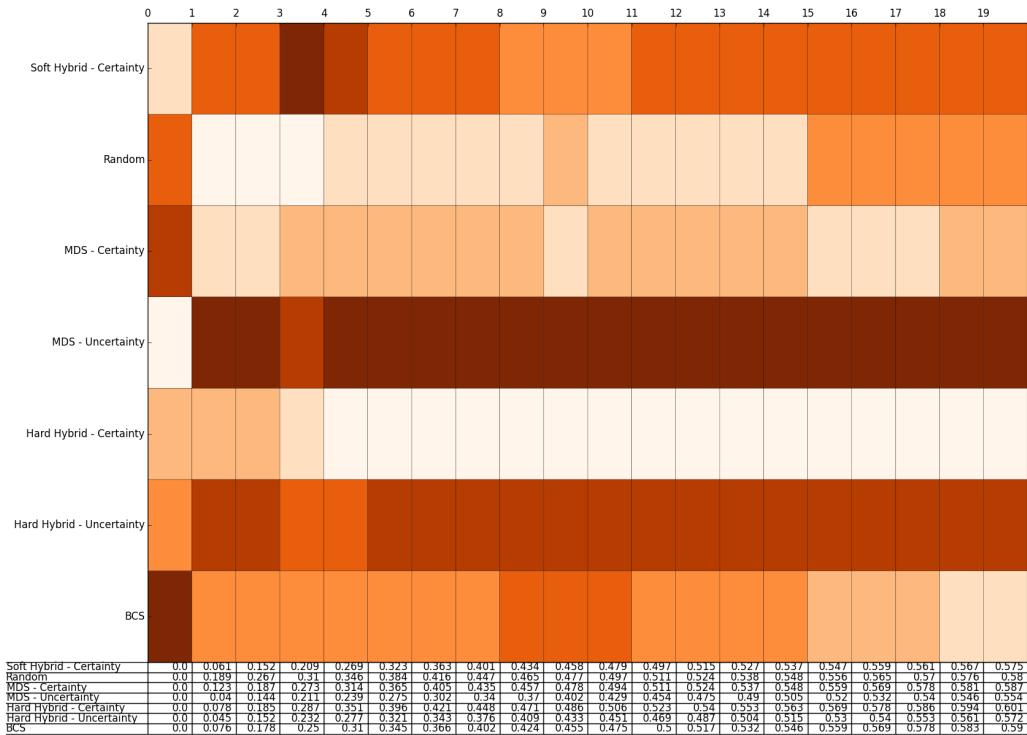
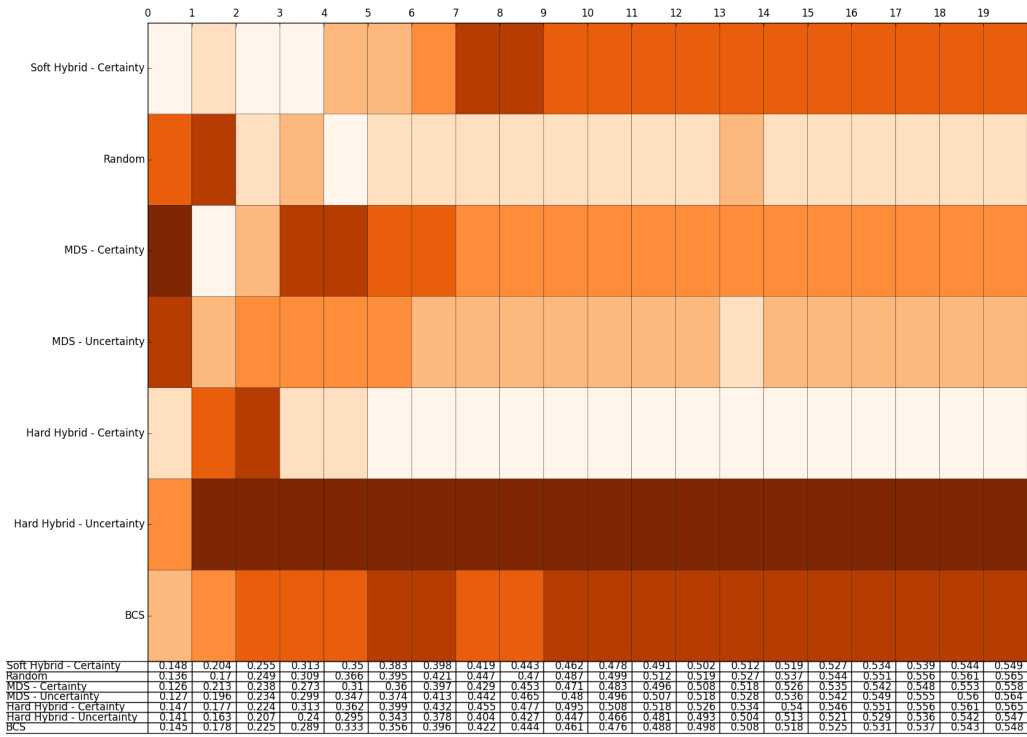
## Speaker score using Speaker annotation is given in orange heat maps.

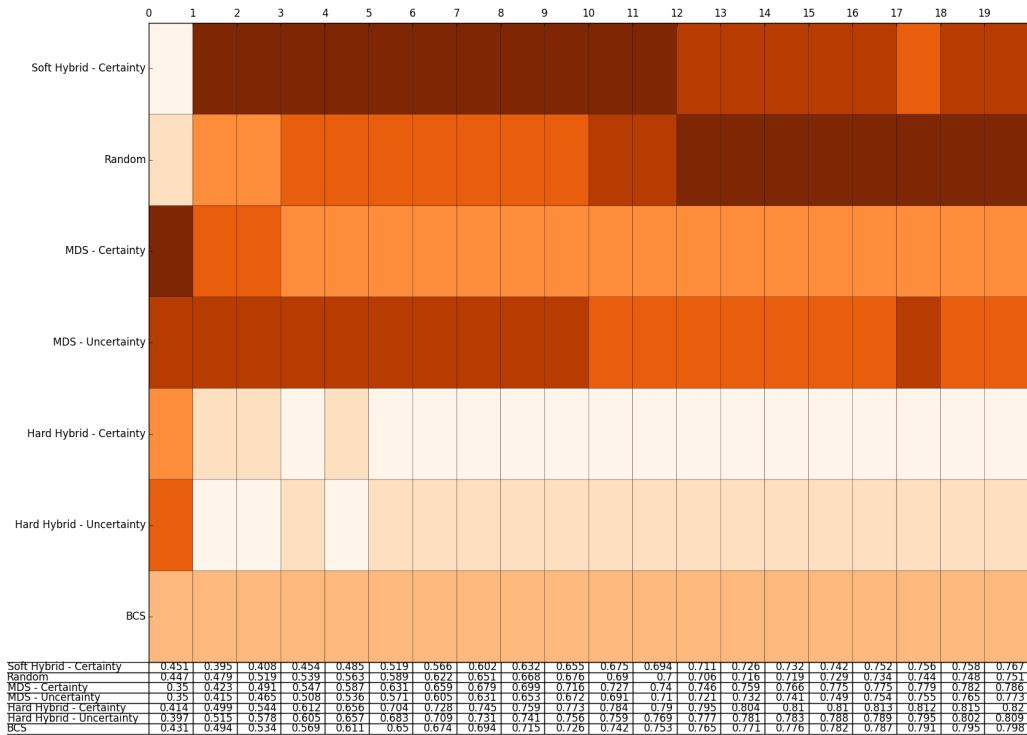
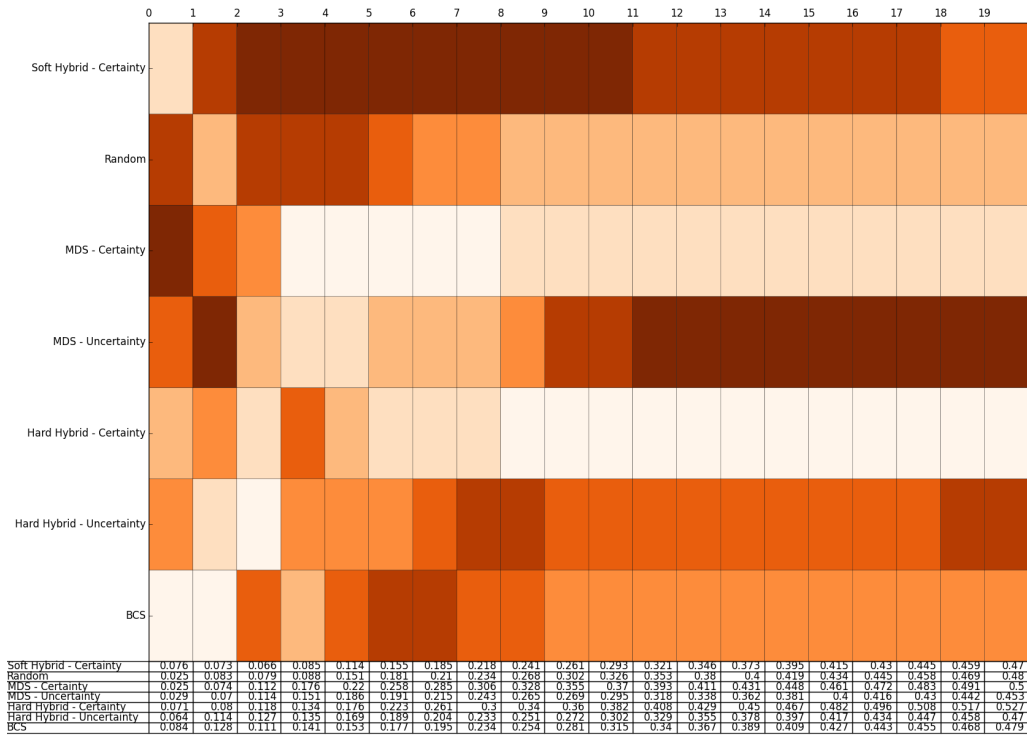




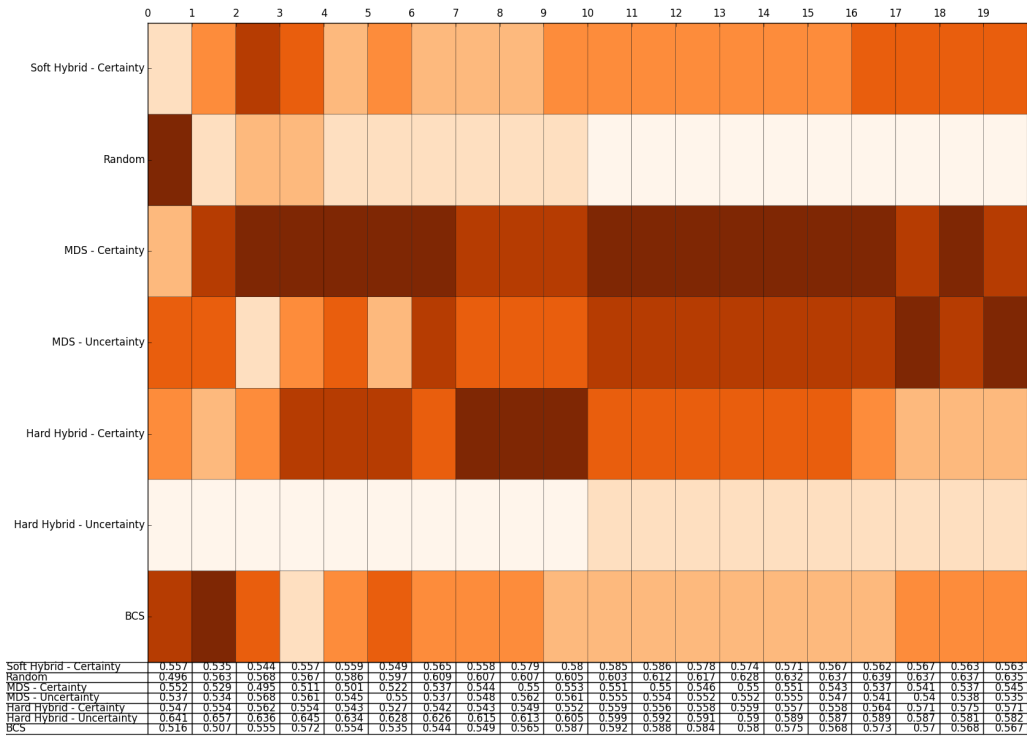
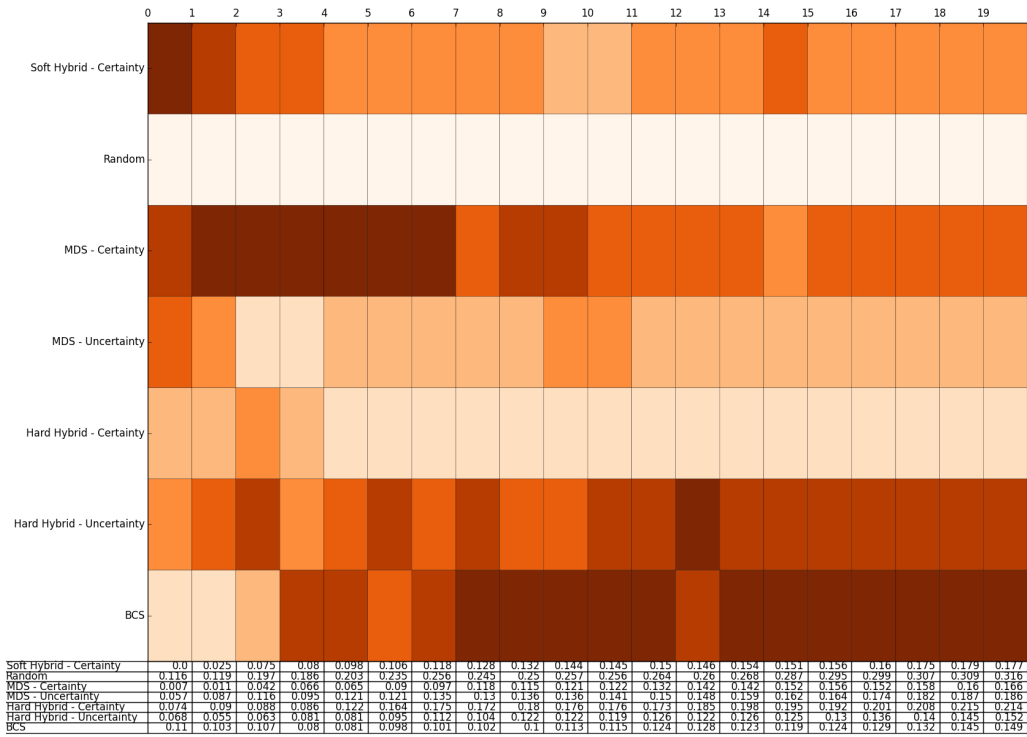


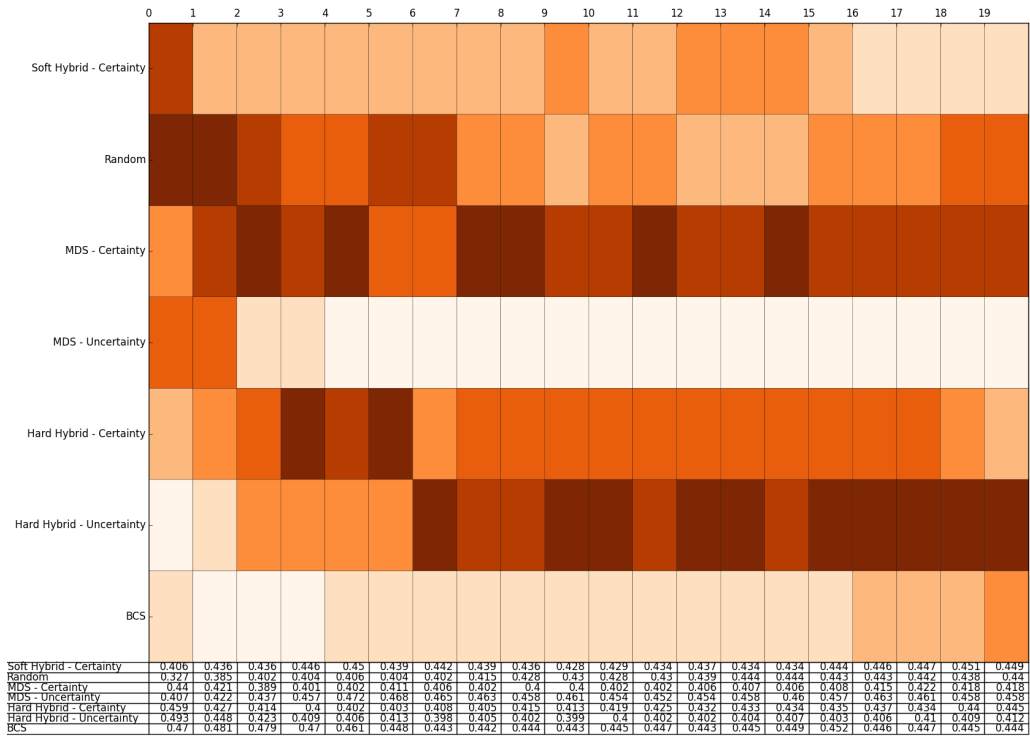
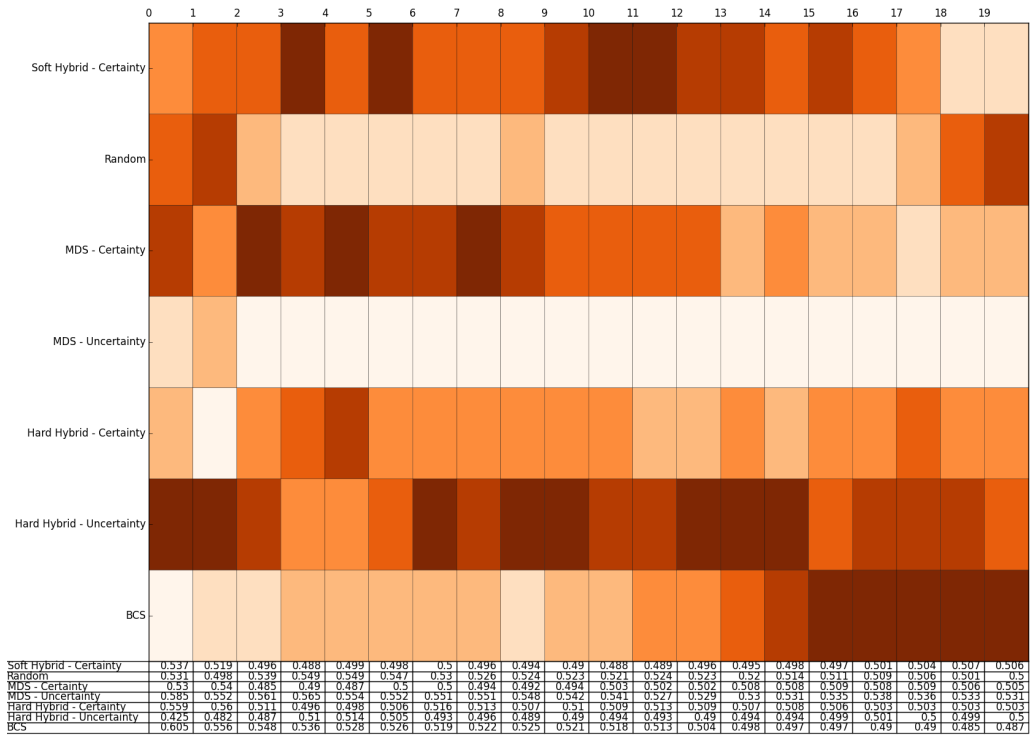


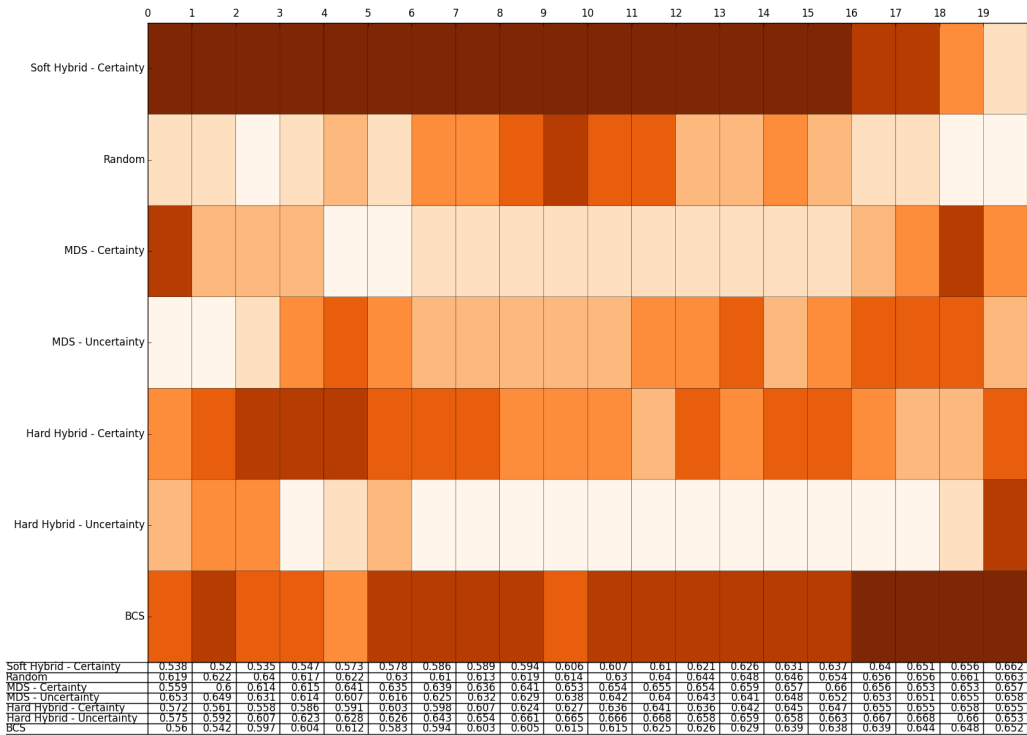
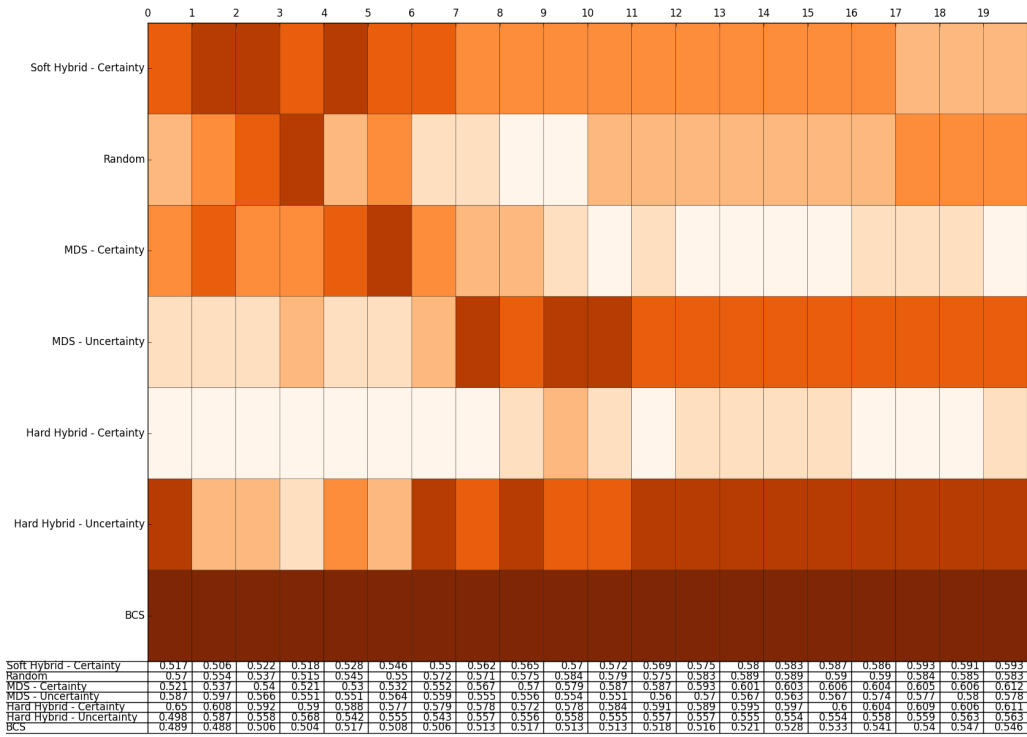


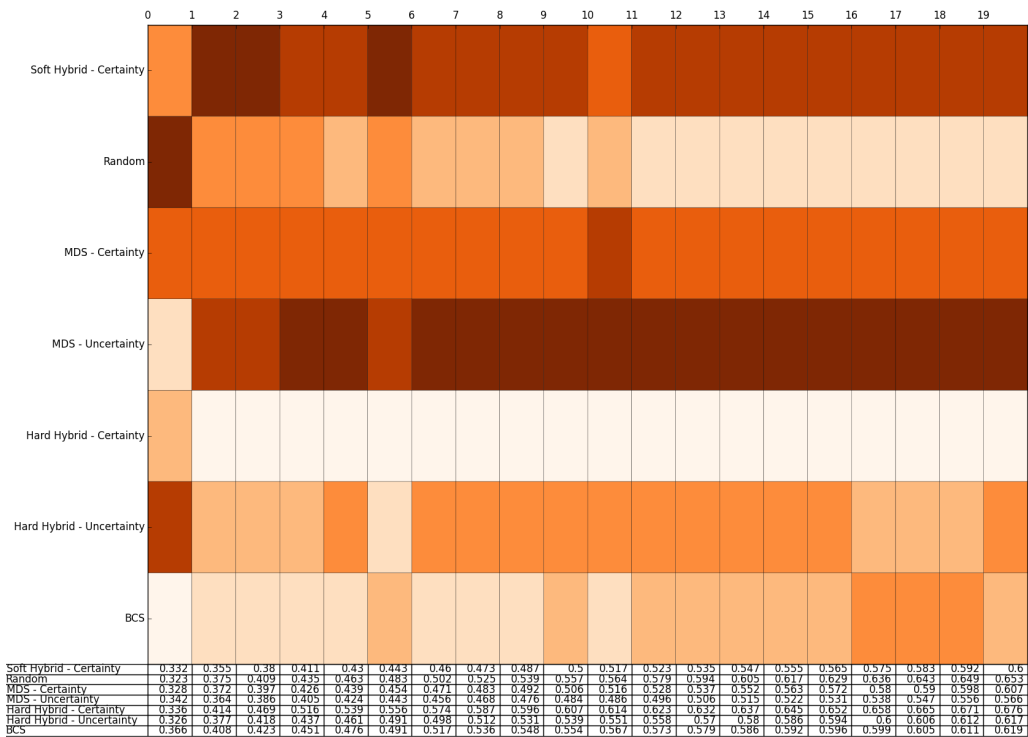
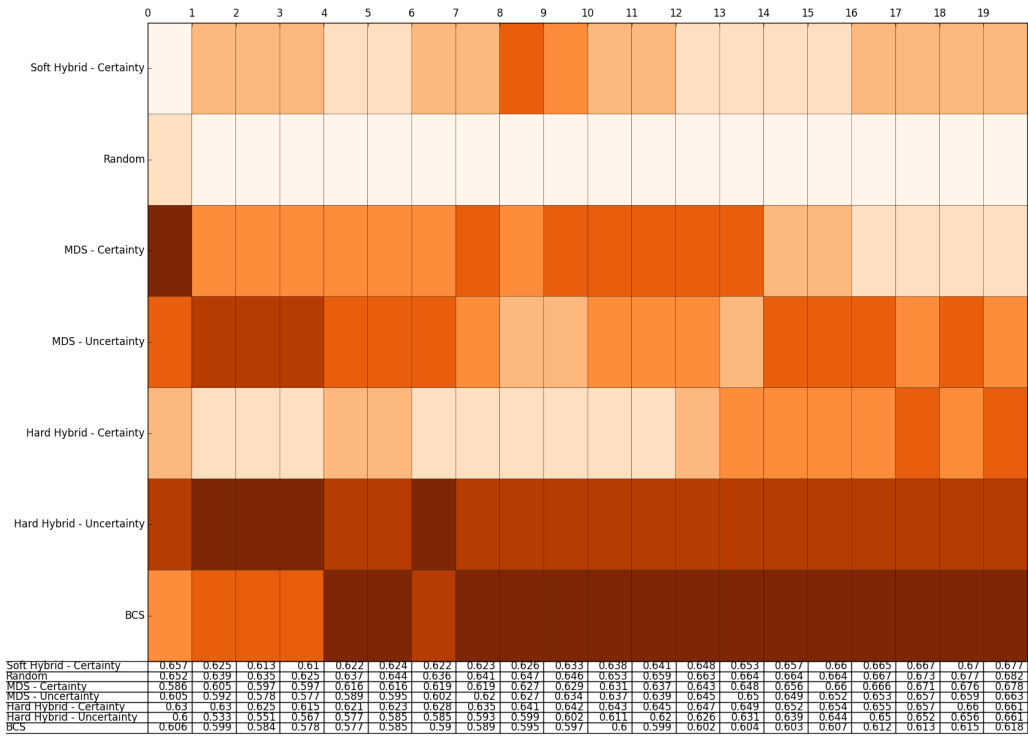


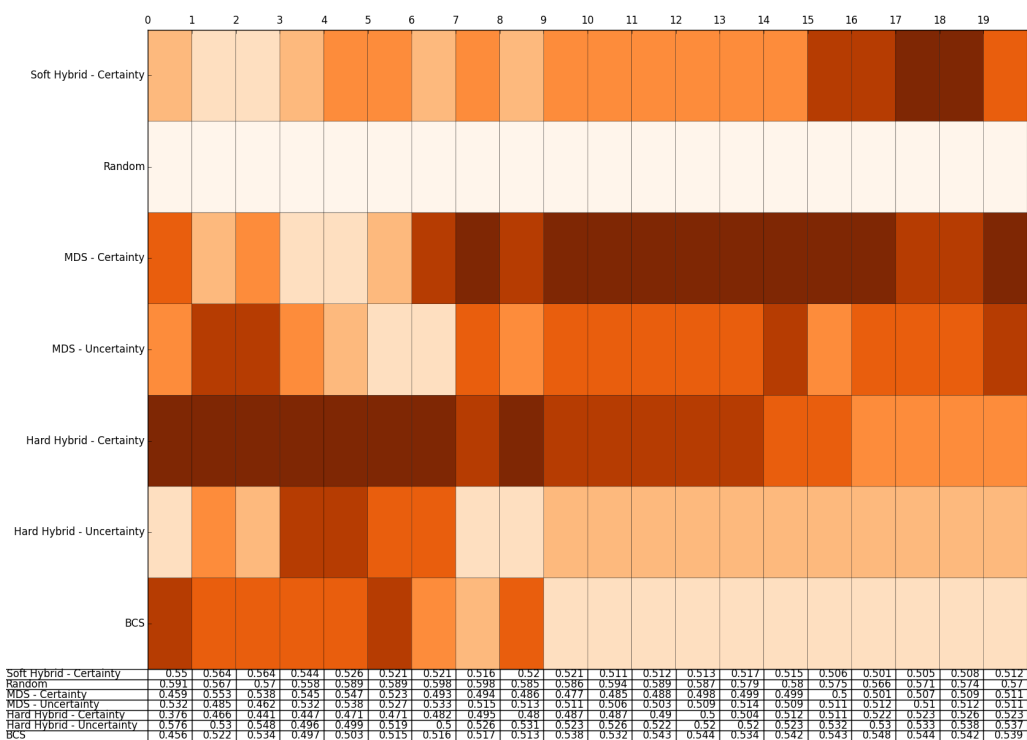
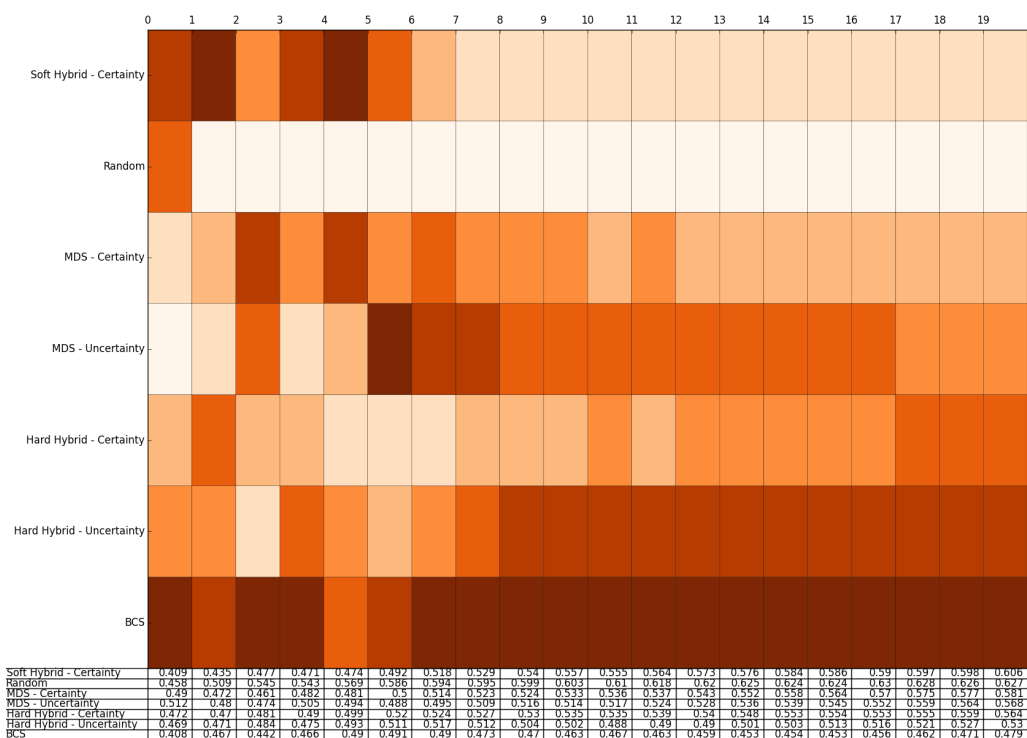


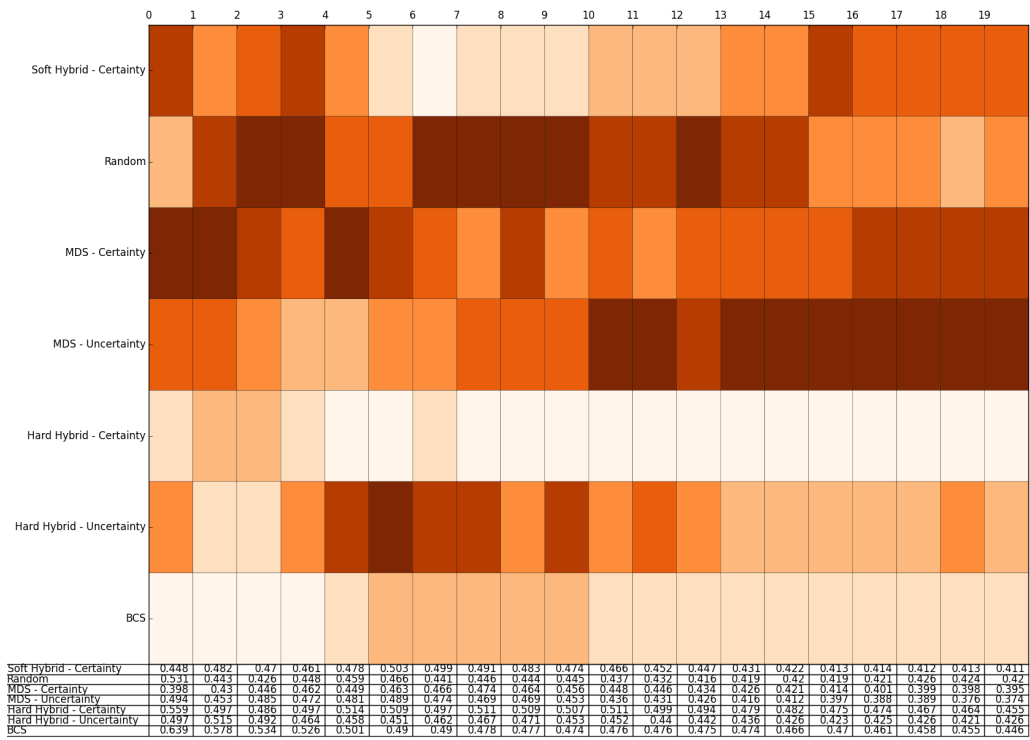
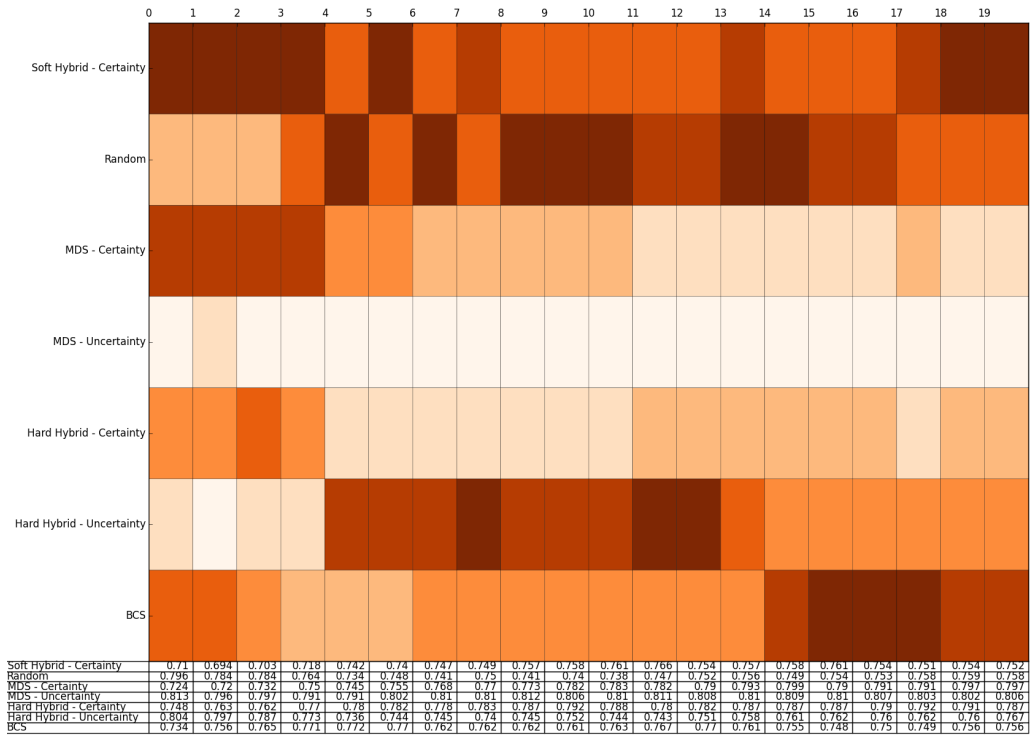


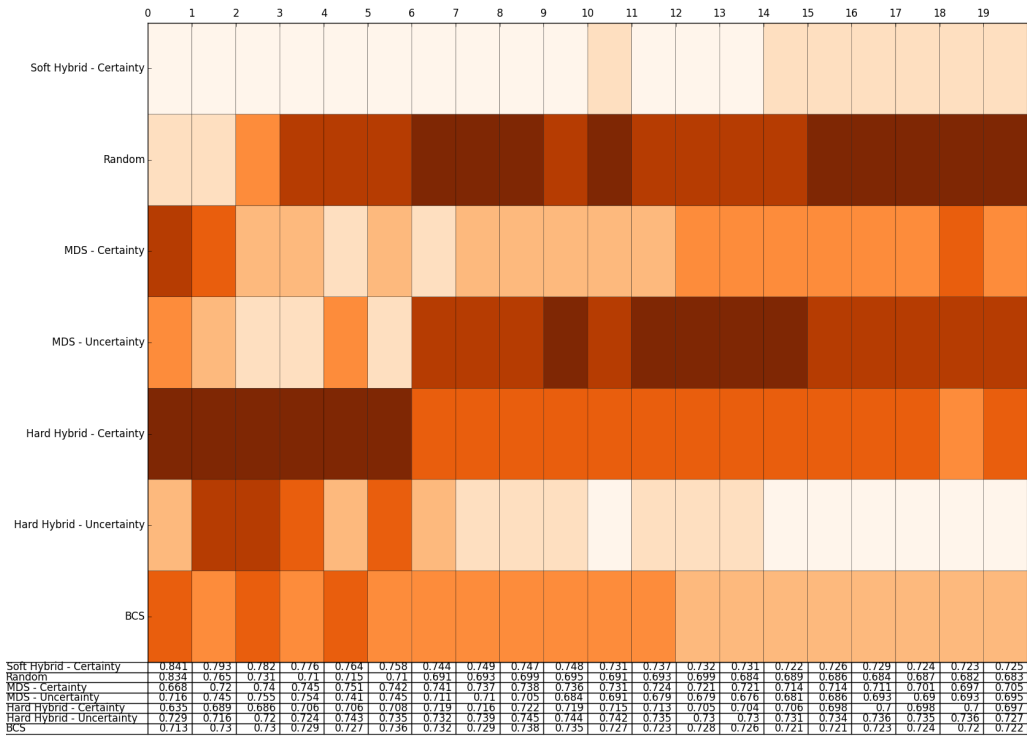
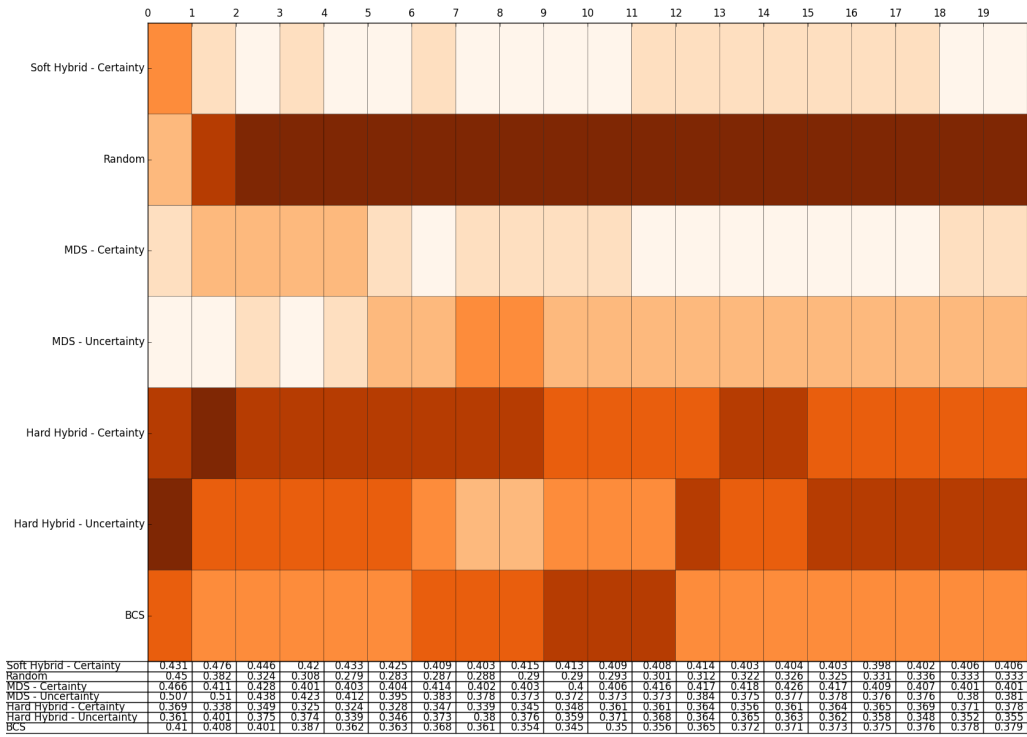


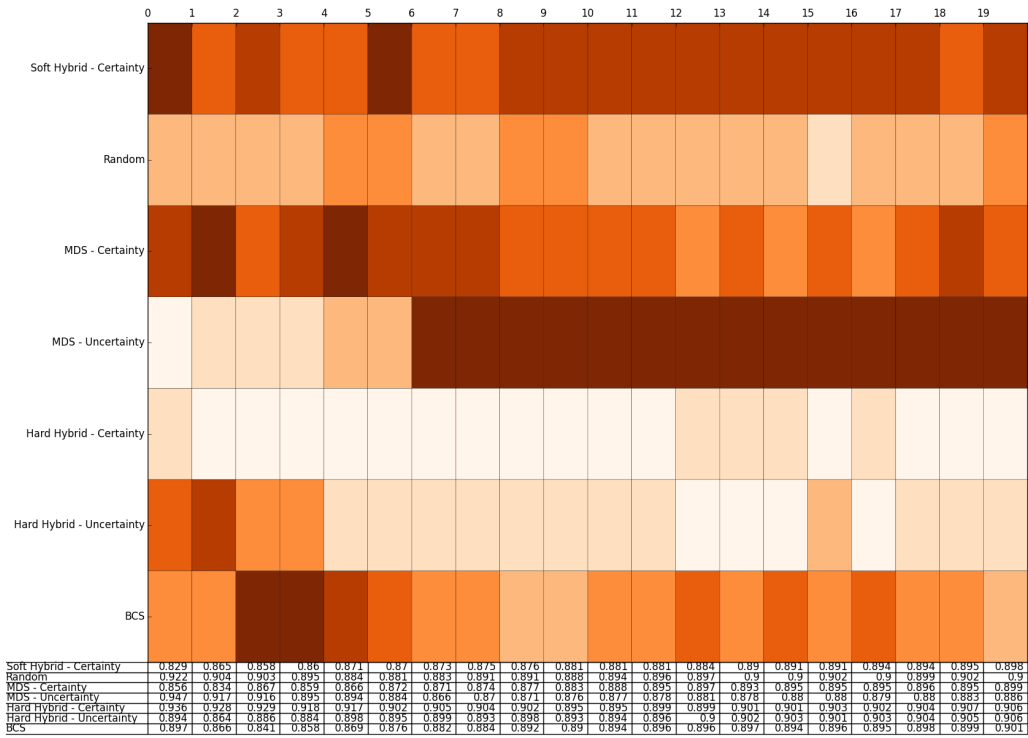
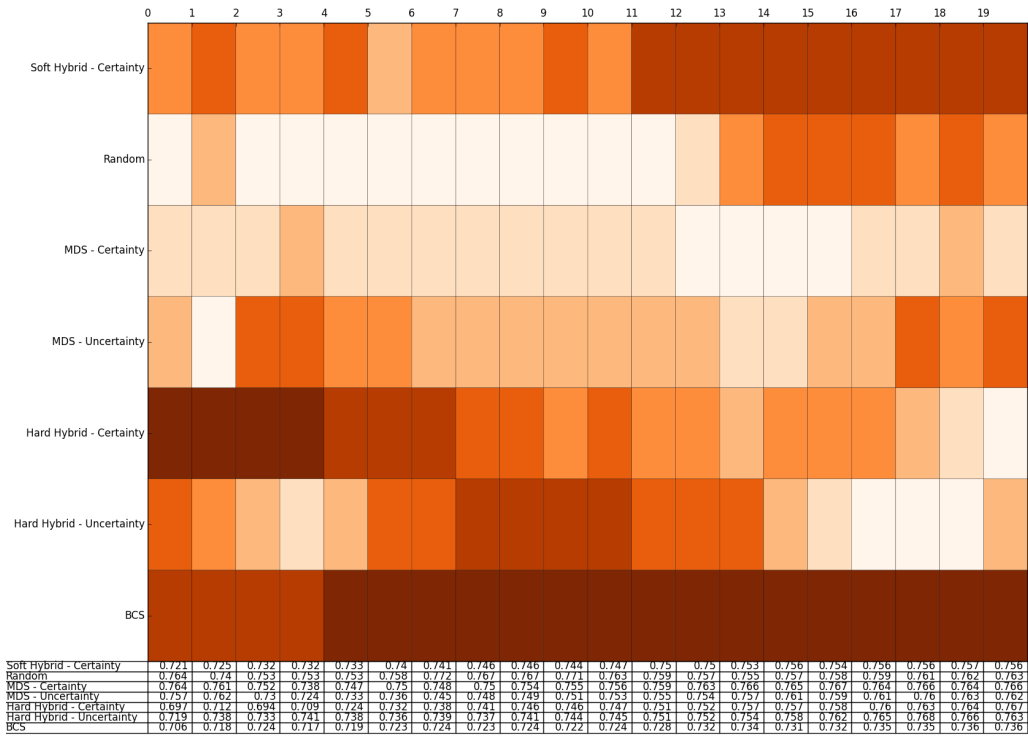




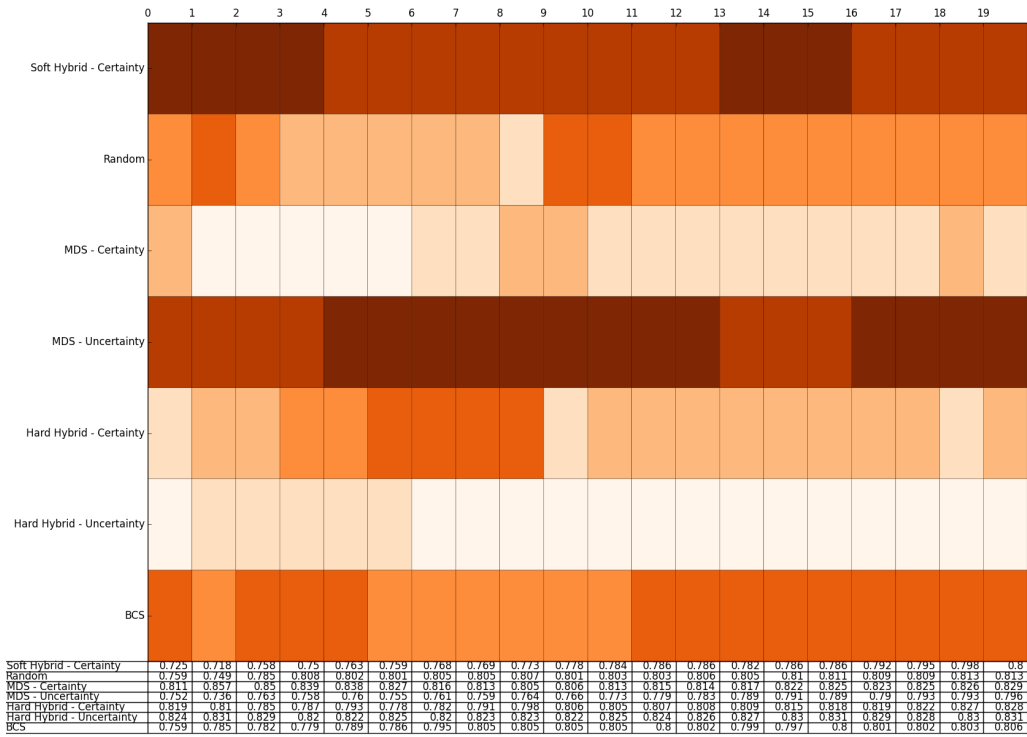














## CURRICULUM VITAE

**Name Surname:** Emre Demir

**Place and Date of Birth:** Zile, 02.08.1985

**Adress:** Tübitak Gebze Yerleşkesi, MAM Lojmanları, 31.Blok, D:5, Gebze/Kocaeli

**E-Mail:** ed.emre@gmail.com

**B.Sc.:** Computer Engineering, Ege University, 2009

### Professional Experience and Rewards:

- 2014 April - Present, Turkcell Teknoloji, R&D department Product and Services division, as Developer.
  - Turkcell Self-Service(ongoing project), sending bulk SMS project, (Developed some parts of project in Java and leading the project).
  - Lawful interception module developed for Turkcell Multi-Media Message Service(MMS), (Java Web Services, Java Crypto Extension).
- 2011 May - 2014 April, The Scientific and Technological Research Council of Turkey, Institute of Cryptology and Electronics - TUBITAK UEKAE, Researcher & Software Engineer.
  - Worked on Electronic Key Management System, (Java Crypto Architecture, Java Crypto Extension, PKI, PKCS 11, Basic Cryptography).
  - Development of Java crypto library, national or NATO key generation algorithms, parity check algorithms etc. implemented (Java Crypto Architecture, Java Crypto Extension).
  - IFFMOD-5 key capability to Electronic Key Management System of Turkish Army.
- 2009 August - 2011 May, Defne Bilgi Islem, Istanbul Office at ITU Teknokent, Software Engineer in R&D Department.
  - Multi-Thread SMPP Server for NCell: High performance SMPP server developed. (Java, SMPP Stack, Sockets, Sonic MQ).
  - Queued Multi-Threaded Diameter Gateway: PayForMe upgrade for Avea (Java, Diameter Stack, Sockets, Linux).



- Multi-Thread Socket Server for IVR and DB Connection Pool, SMPP Client, Web Services, Back Ground Music Project, for Mobilink Pakistan (Java, Tomcat, SMPP Stack, Sockets, Linux).
- Web Services, SMPP Client, Malicious Call Barring Project for Azercell (Java, Tomcat, SMPP Stack, Linux).
- Making Videos From Web Pages, Research Project (FFMpeg API, C, Linux).
- Video Transcoding and Streaming, IP TV Research Project (FFMpeg API, C, Linux).

**List of Publications and Patents:**

**PUBLICATIONS/PRESENTATIONS ON THE THESIS**

- Unsupervised Active Learning for Video Annotation, International Conference of Machine Learning 2015, Active Learning Workshop, Lille - France.