



# **Database Schema**

## **(Corpora and annotation tools)**

Prepared by: Yoong Kuan Goh (Andrew)

Department: Connected Intelligence Center in UTS

## Contents

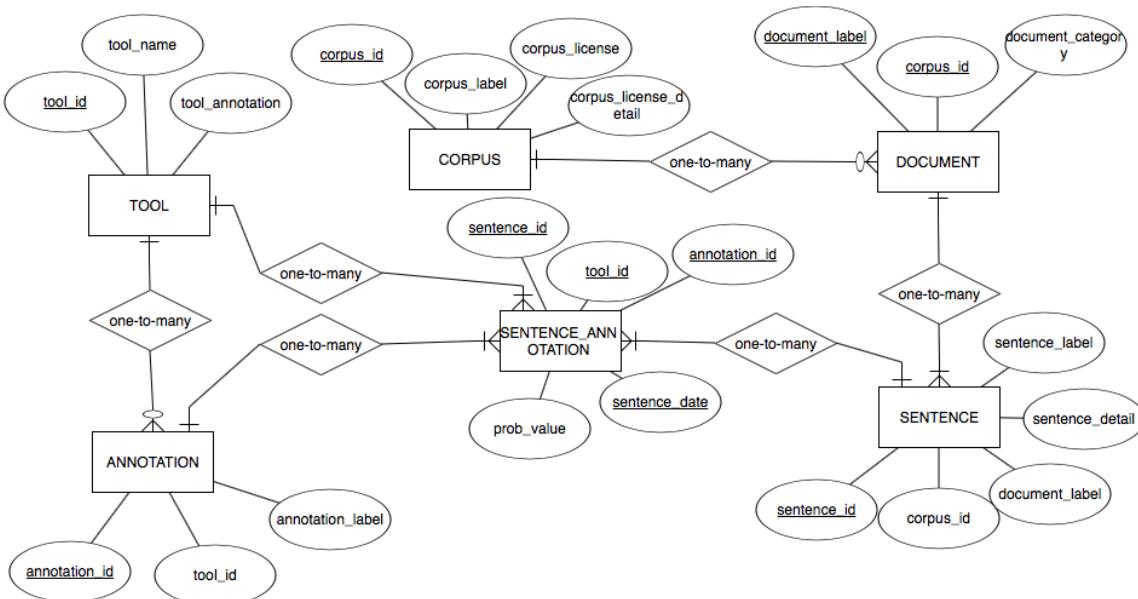
<b>1. Introduction.....</b>	<b>2</b>
<b>2.0 Annotation tools .....</b>	<b>2</b>
<b>3.0 Annotation types .....</b>	<b>3</b>
<b>4.0 Corpus.....</b>	<b>4</b>
<b>5.0 Corpus (Document).....</b>	<b>4</b>
<b>6.0 Sentences.....</b>	<b>6</b>
<b>7.0 Sentence Annotations .....</b>	<b>8</b>
<b>7.1 AntMover annotation distribution .....</b>	<b>10</b>
<b>7.2 AWA annotation distribution .....</b>	<b>10</b>
<b>7.3 AWA3 annotation distribution .....</b>	<b>10</b>
<b>7.4 RWT annotation distribution .....</b>	<b>11</b>
<b>8.0 Association within AWA .....</b>	<b>11</b>
<b>9.0 Association within AWA3 .....</b>	<b>13</b>
<b>10.0 AWA and AWA3 Associations .....</b>	<b>15</b>
<b>10.1 AWA vs AWA3 .....</b>	<b>15</b>
<b>10.2 AWA3 vs AWA .....</b>	<b>16</b>
<b>11. AWA3 and RWT associations .....</b>	<b>19</b>
<b>11.1 AWA3 VS RWT .....</b>	<b>19</b>
<b>11.2 RWT VS AWA3 .....</b>	<b>20</b>
<b>12. AWA3 and RWT associations (Ignoring sentences that not annotated) .....</b>	<b>22</b>
<b>12.1 AWA3 VS RWT .....</b>	<b>22</b>
<b>12.2 RWT VS AWA3 .....</b>	<b>23</b>
<b>13. AWA3 and AntMover associations .....</b>	<b>25</b>
<b>13.1 AWA3 VS AntMover .....</b>	<b>25</b>
<b>13.2 AntMover VS AWA3 .....</b>	<b>26</b>
<b>14. Conclusion .....</b>	<b>28</b>

## 1. Introduction

This report describes the structure of a database scheme that been setup for the annotation analysis project involving four different annotation tools across three different corpora. Each corpus is broken down to document and sentence level. A MySQL database is used to store the annotation results.

Database name: AWA

ERD diagram:



## 2.0 Annotation tools

The annotation tools are AntMover, AWA, AWA3 and RWT. The information of these annotation tools is being stored in a table called ‘TOOL’.

Field	Type	Null	Key	Default	Extra
tool_id	int(11)	NO	PRI	NULL	auto_increment
tool_name	varchar(50)	YES		NULL	
tool_annotation	varchar(50)	YES		NULL	

Fields:

- a) tool\_id → store the id of each tool
- b) tool\_name → store the name of each tool
- c) tool\_annotation → a short description of each tool

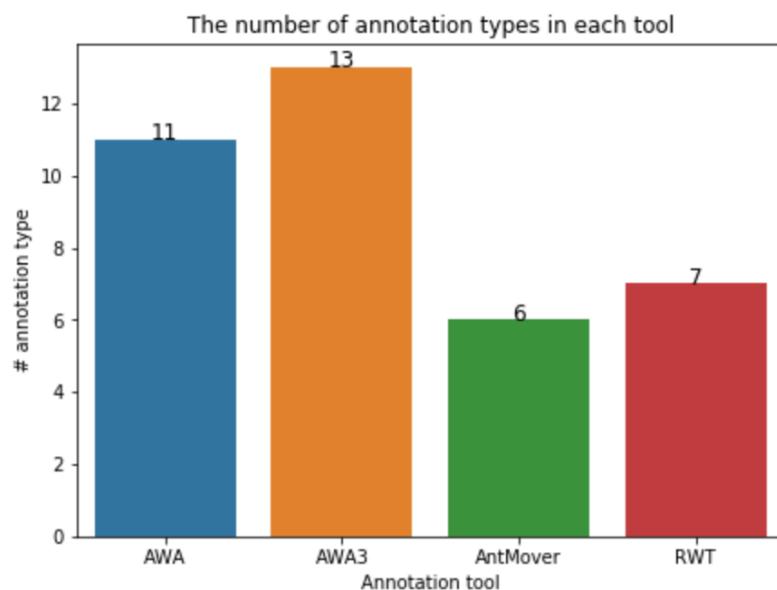
### 3.0 Annotation types

Each annotation tool has specific annotation scheme. The number of annotation type in each scheme may be different. These data are stored in a table called ‘ANNOTATION’ and this table is the child table of ‘TOOL’ table.

Field	Type	Null	Key	Default	Extra
annotation_id	int(11)	NO	PRI	NULL	auto_increment
annotation_label	varchar(50)	YES		NULL	
tool_id	int(11)	YES	MUL	NULL	

Fields:

- a) annotation\_id → store the id of each annotation type
- b) annotation\_label → store the name of each annotation type
- c) tool\_id → store the referencing tool\_id from the ‘TOOL’ table



## 4.0 Corpus

The corpora used as the raw data of the annotation analysis are Elsevier (OA-STM), BAWE and PMC. The table that stored the type of corpus is ‘CORPUS’.

Field	Type	Null	Key	Default	Extra
► corpus_id	int(11)	NO	PRI	NULL	auto_increment
corpus_label	varchar(50)	YES		NULL	
corpus_license	tinyint(1)	YES		NULL	
corpus_license_detail	varchar(250)	YES		NULL	

Fields:

- a) corpus\_id → store the id of each corpus
- b) corpus\_label → store the name of each corpus
- c) corpus\_license → state whether licensing is required for the corpus. 0 = no and 1 = yes
- d) corpus\_license\_detail → describe about the license.

## 5.0 Corpus (Document)

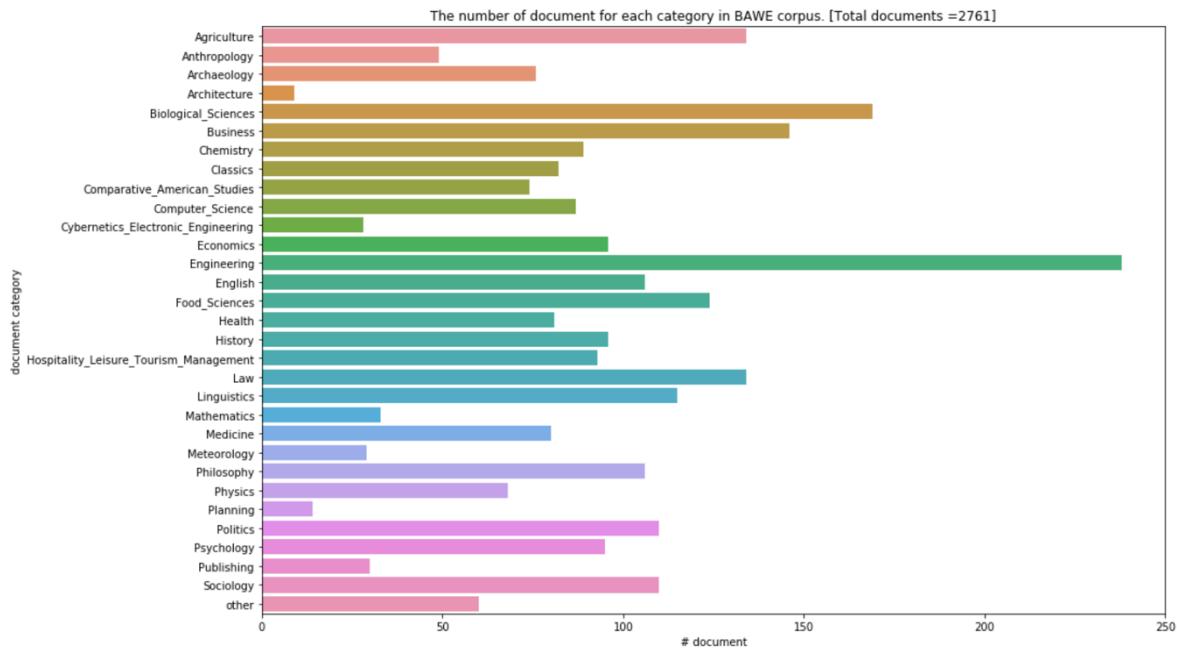
Each corpus contains a lot of files. Each file is called as a document. The document details are stored in a table called ‘DOCUMENT’. This table is a child table of ‘CORPUS’ table.

Field	Type	Null	Key	Default	Extra
► document_label	varchar(100)	NO	PRI	NULL	
document_category	varchar(50)	YES		NULL	
corpus_id	int(11)	NO	PRI	NULL	

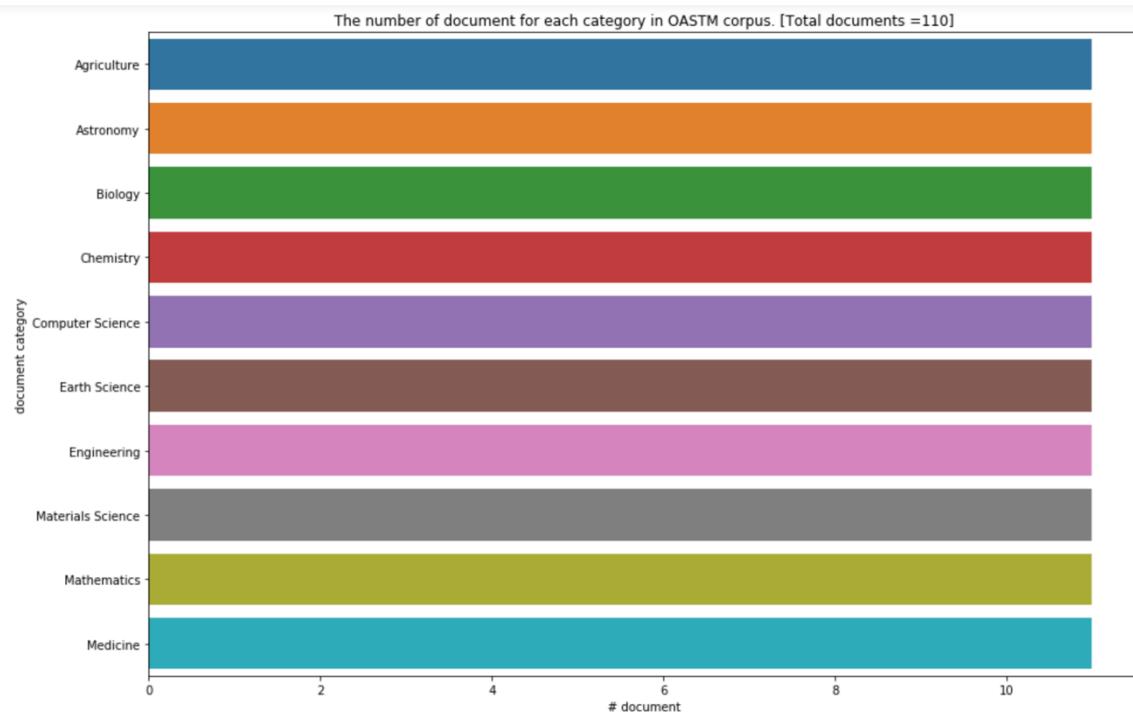
Fields:

- a) document\_label → store the name of the file/document
- b) document\_category → store the category of the document
- c) corpus\_id → store the referencing corpus\_id from the ‘TOOL’ table

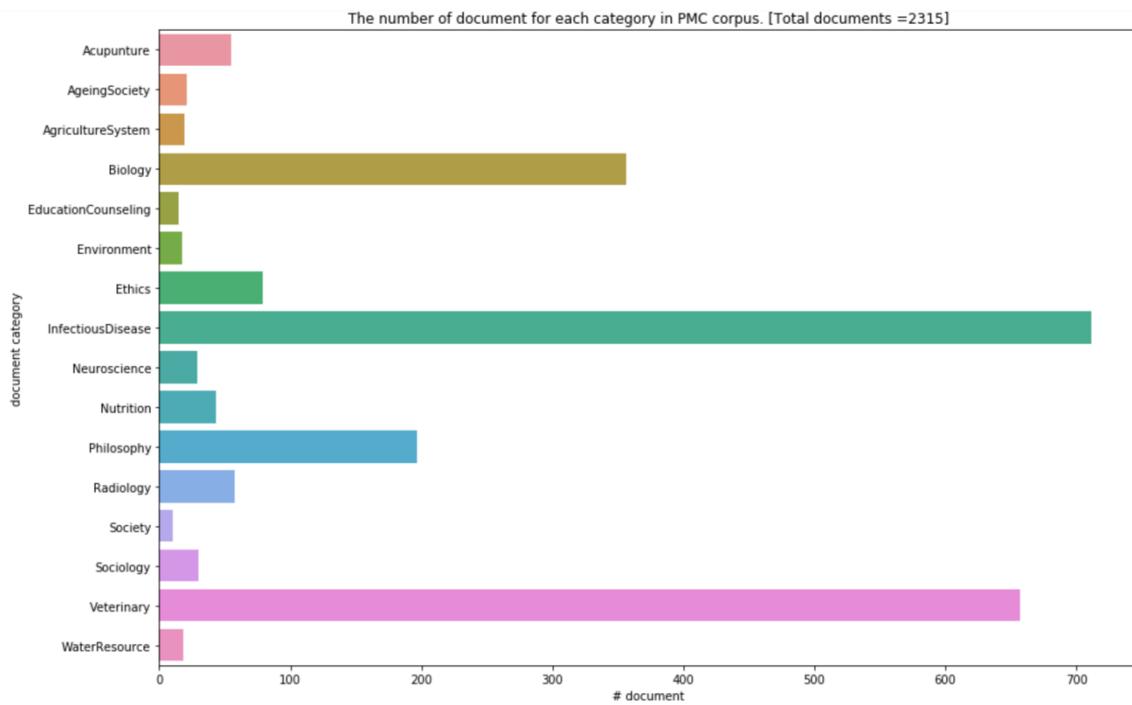
### a) BAWE corpus



### b) Elsevier/OA-STM corpus



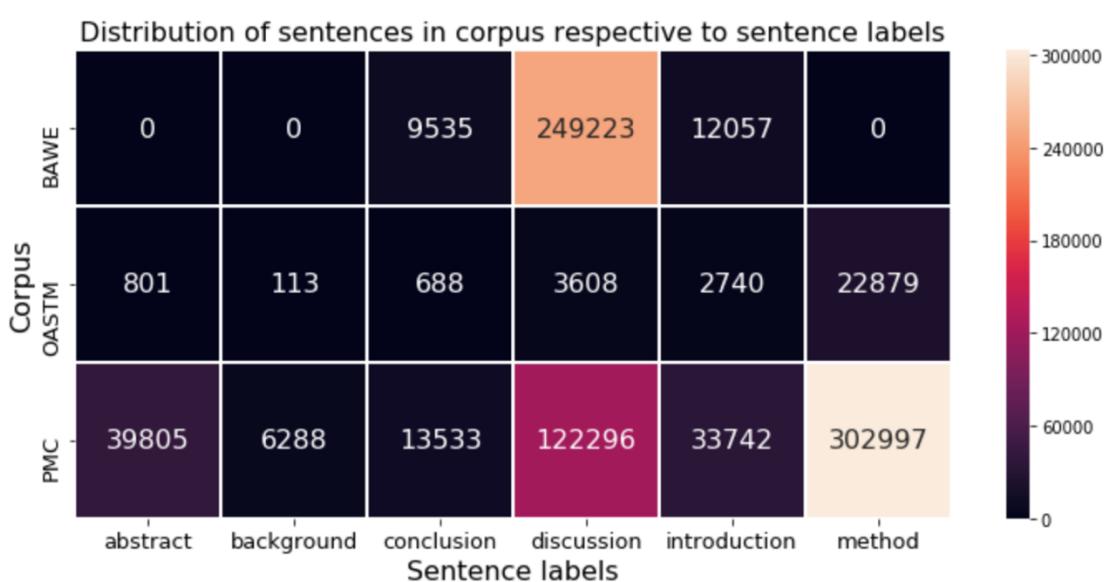
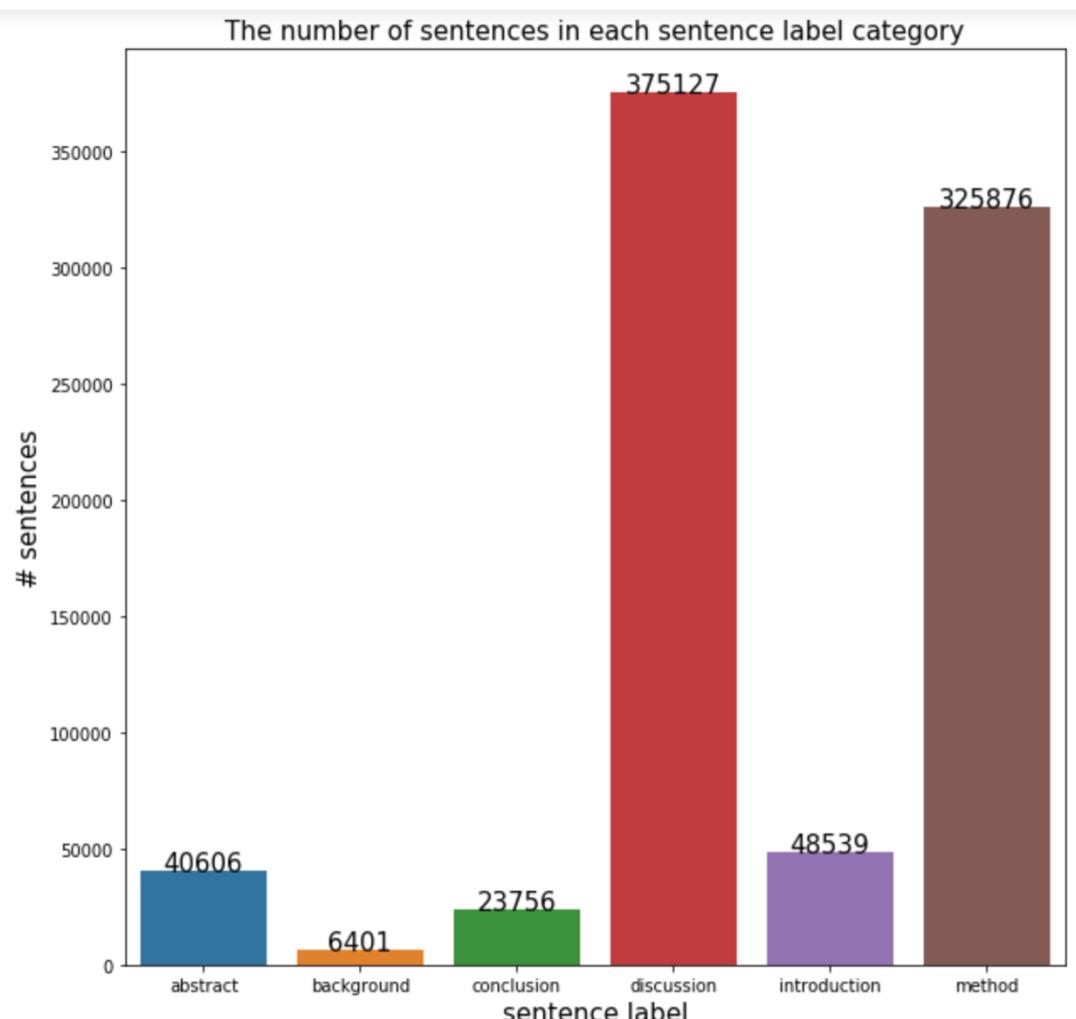
### c) PMC corpus



## 6.0 Sentences

Each document is broken down into sentences. Each sentence is stored in a table called ‘SENTENCE’. The document labels of the sentences are stored together in the same table for easy referencing. Each sentence can be categorized into either ‘abstract’, ‘introduction’, ‘background’, ‘discussion’, ‘method’ or ‘conclusion’ as sentence label.

Field	Type	Null	Key	Default	Extra
▶ sentence_id	bigint(20)	NO	PRI	NULL	auto_increment
sentence_detail	varchar(2500)	YES	MUL	NULL	
sentence_label	varchar(350)	YES		NULL	
document_label	varchar(100)	YES	MUL	NULL	
corpus_id	int(11)	YES		NULL	



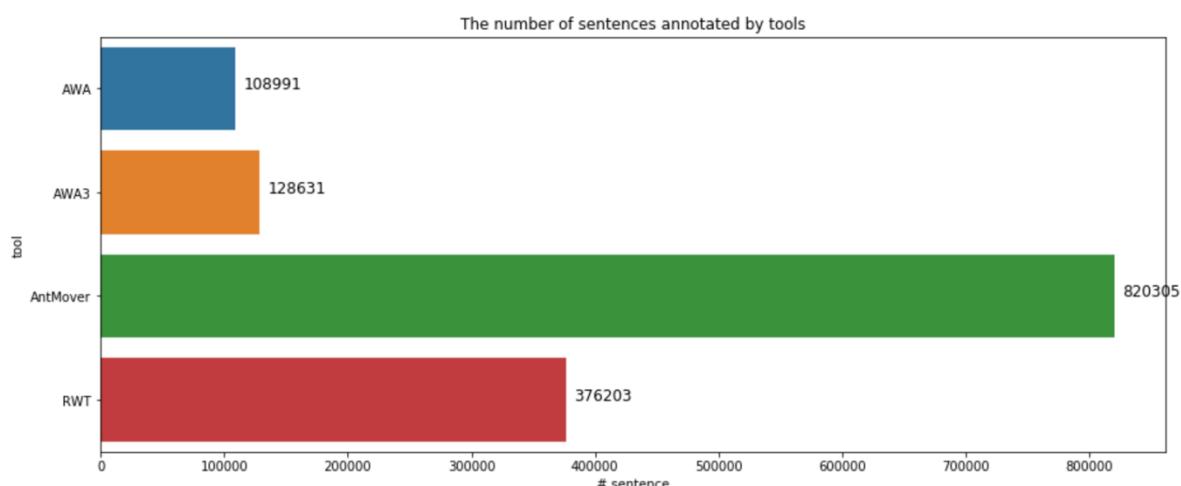
Total number of sentences: **820,305**

Total number of distinct sentences: **629,275**

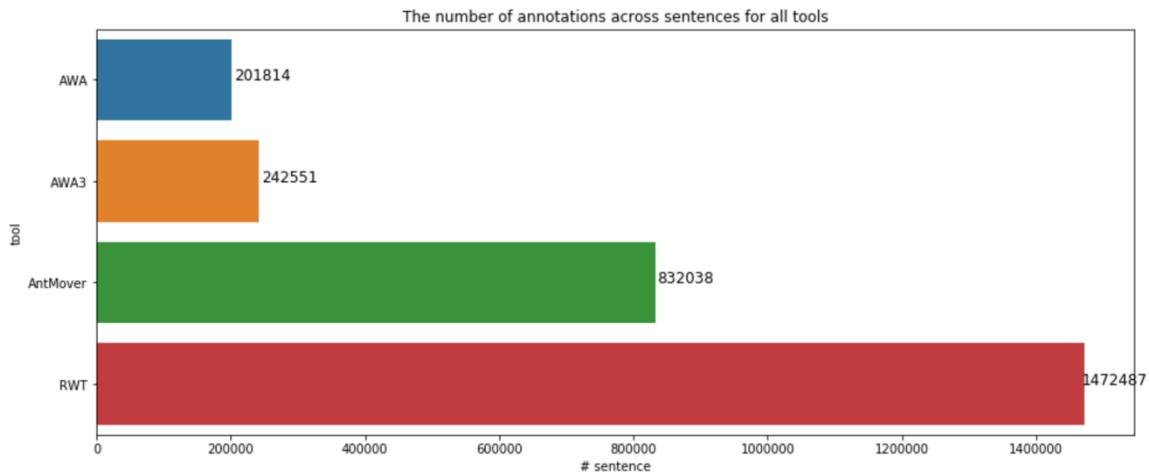
## 7.0 Sentence Annotations

Each sentence in ‘SENTENCE’ table will be annotated with all the tools from the ‘TOOL’ table. The annotated sentences will be stored in ‘SENTENCE\_ANNOTATION’ table. Since RWT annotation involves the probabilistic values for each of its annotation, a column called ‘prob\_value’ is created to store the probabilities. Sentences annotated by other tools apart from RWT will always have probability value as 1.

Field	Type	Null	Key	Default	Extra
▶ sentence_id	bigint(20)	NO	PRI	NULL	
tool_id	int(11)	NO	PRI	NULL	
sentence_date	date	NO	PRI	NULL	
annotation_id	int(11)	NO	PRI	NULL	
prob_value	float	YES		0	



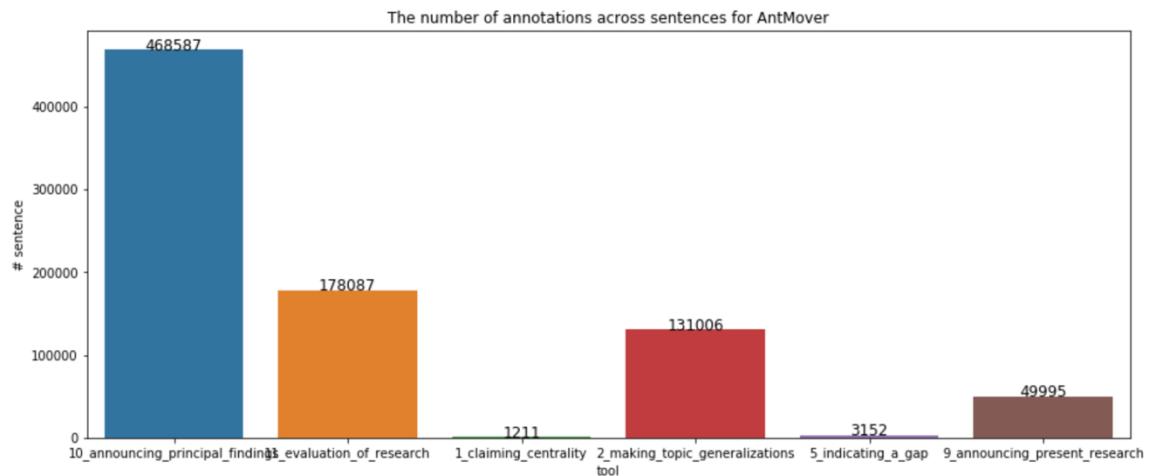
From the above graph, not all tools will annotate every sentence. Only AntMover tool will annotate every sentence. The tool that annotate the least sentences is AWA followed by AWA3 and RWT. The difference between AWA and AWA3 is 18.02% (19640 sentences). In total **820,305** sentences, **AWA3 only annotate 15.68% of the sentences**. This shows that there are still plenty of improvement or enhancement can be made to AWA3 so that it can recognize and annotate more sentences.



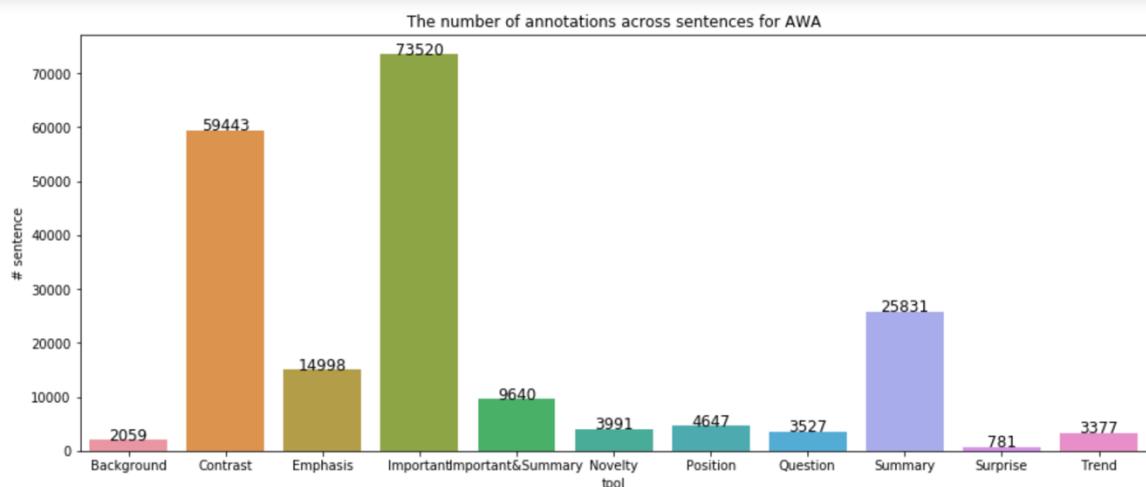
The above graph shows the total number of annotations in all the sentences. The RWT has the highest number of annotations because each sentence is annotated and given different probabilities values for each annotation. As an example, sentence A would be annotated and be given 7 probabilities value that show how likely this sentence being classify in each RWT annotations. So, a sentence will be annotated with 7 different annotation with different probability value.

Meanwhile, even though AntMover is supposed to annotated every annotation with only one annotation but surprisingly the number of annotations is higher than the total number of sentences. This is because, AntMover annotated a same sentence but in different files with different annotations. As an example, a sentence “Click here for additional data file” from file A had been annotated with “making\_topic\_generalizations” annotation but being annotated with “announcing\_present\_research” annotation at file B. Therefore, a same sentence had been considered annotated twice in this case.

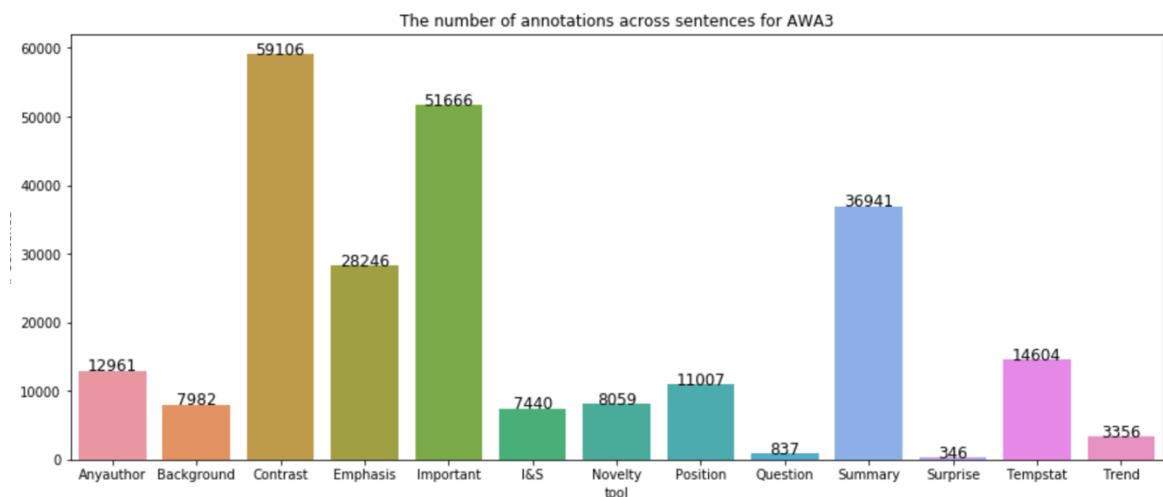
## 7.1 AntMover annotation distribution



## 7.2 AWA annotation distribution

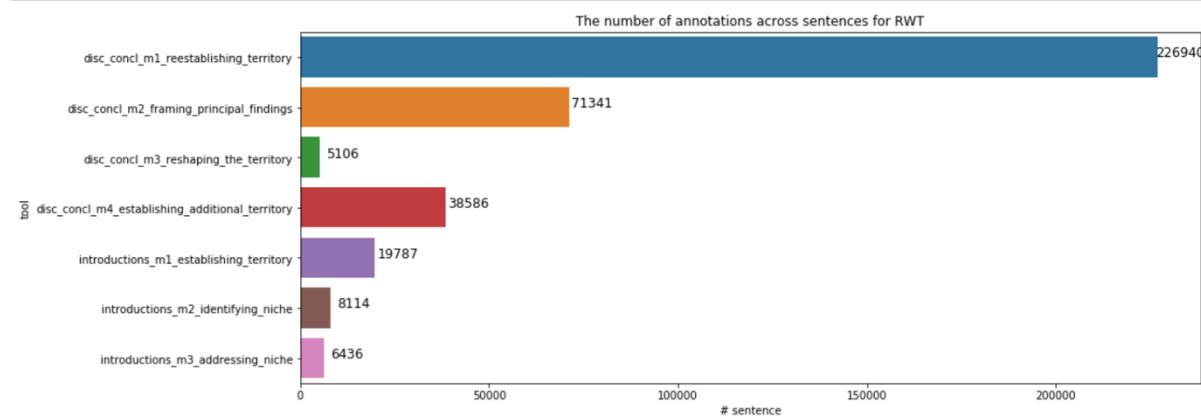


## 7.3 AWA3 annotation distribution



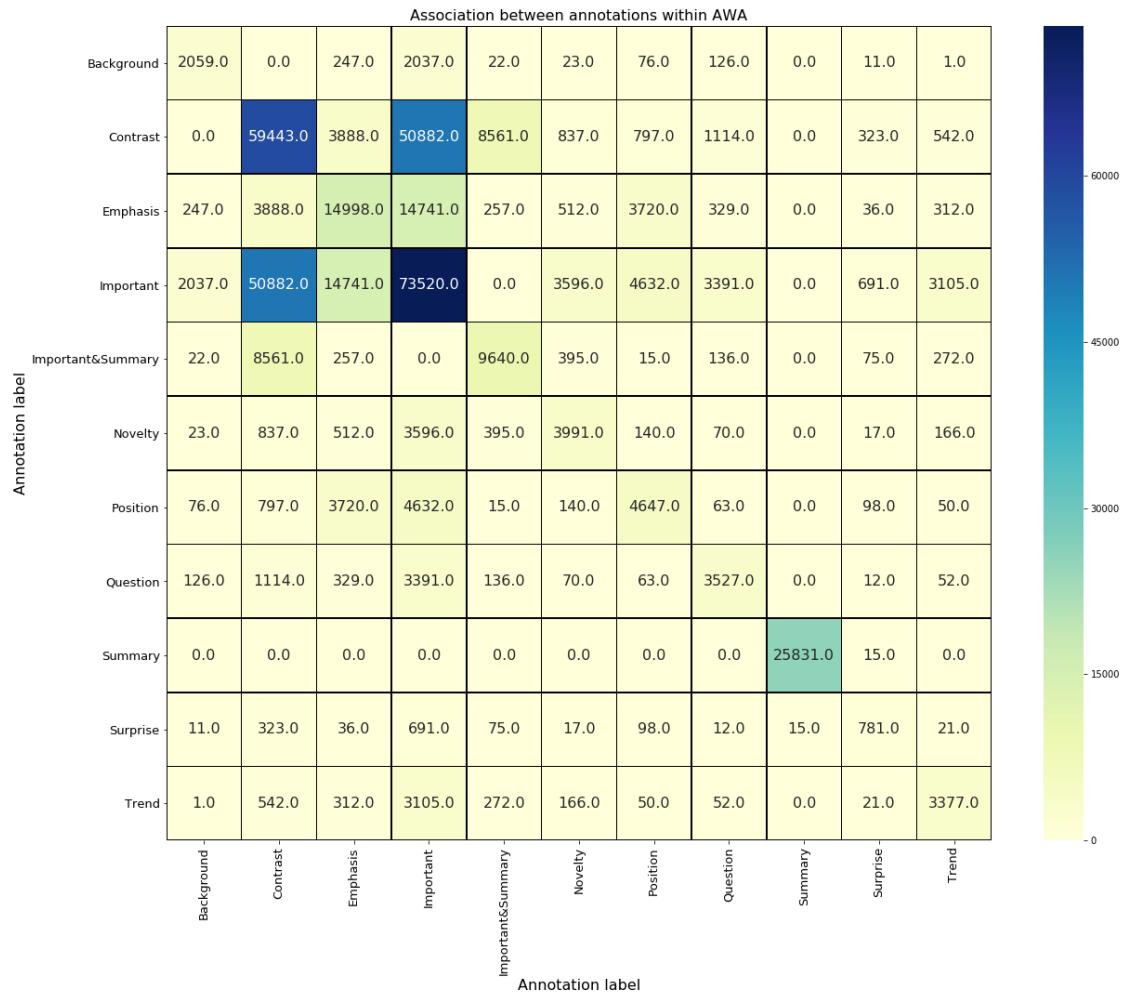
## 7.4 RWT annotation distribution

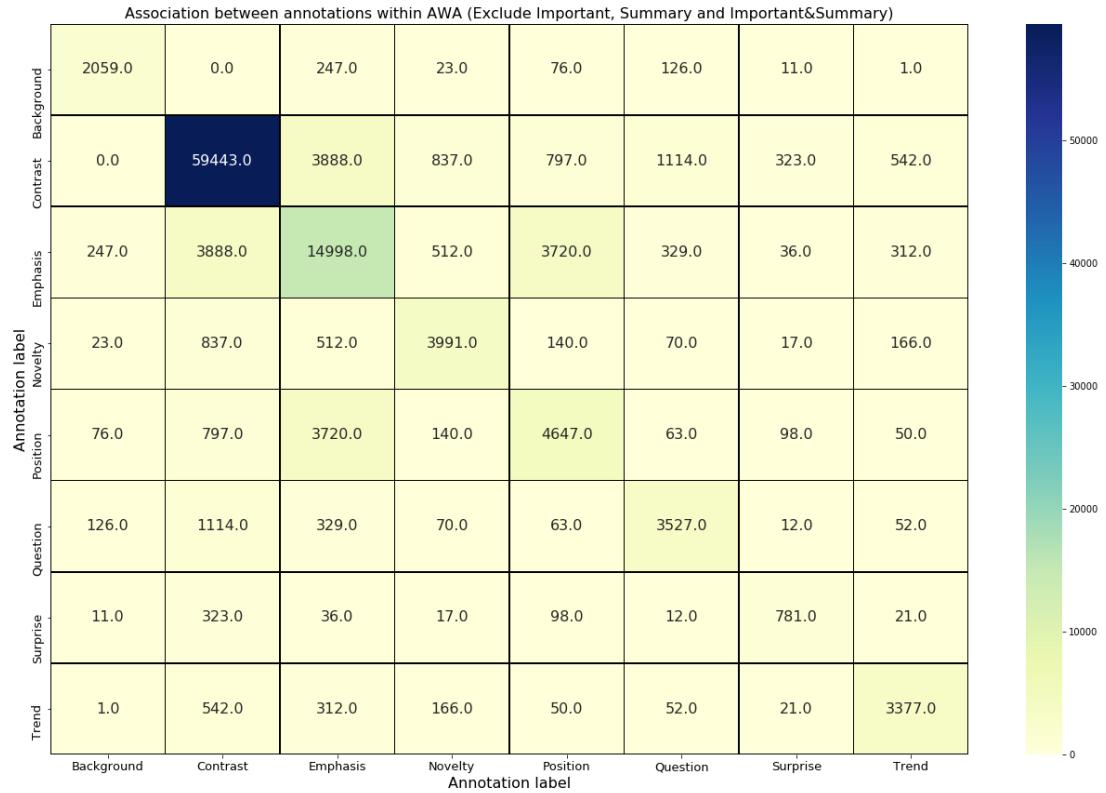
As each sentence would be given probability value for each annotation, only the record with the highest probability for each sentence will be taken in this distribution.



## 8.0 Association within AWA

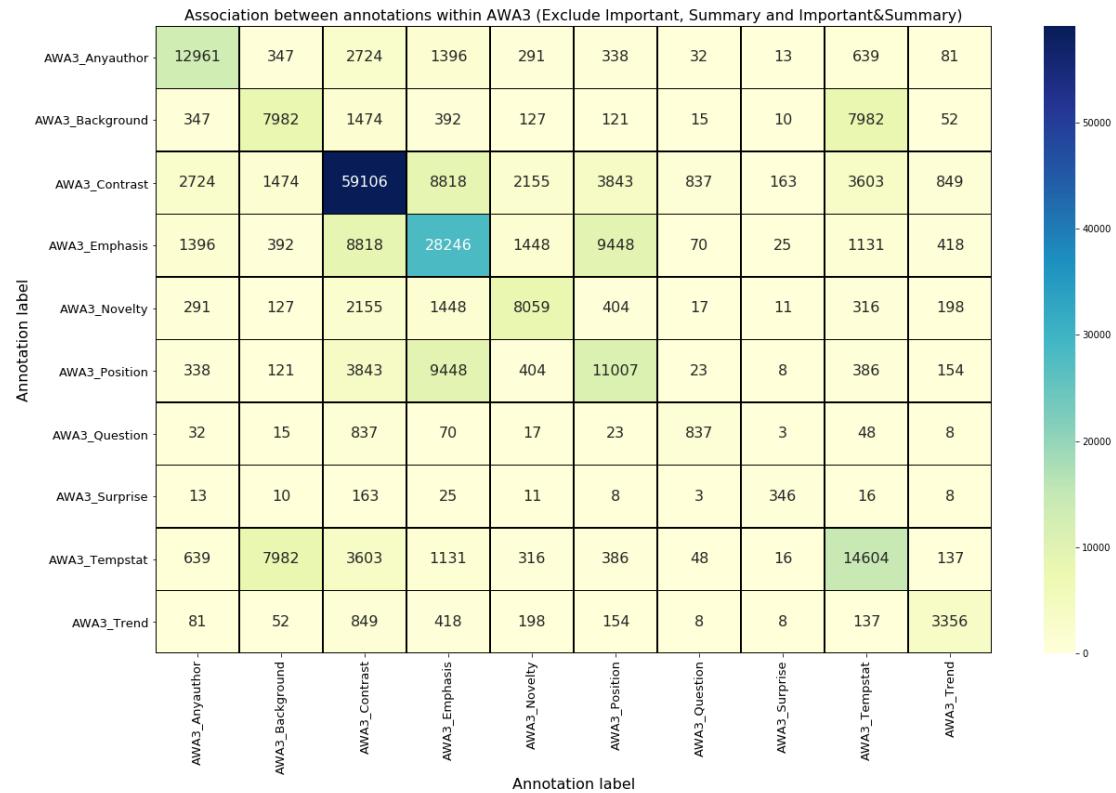
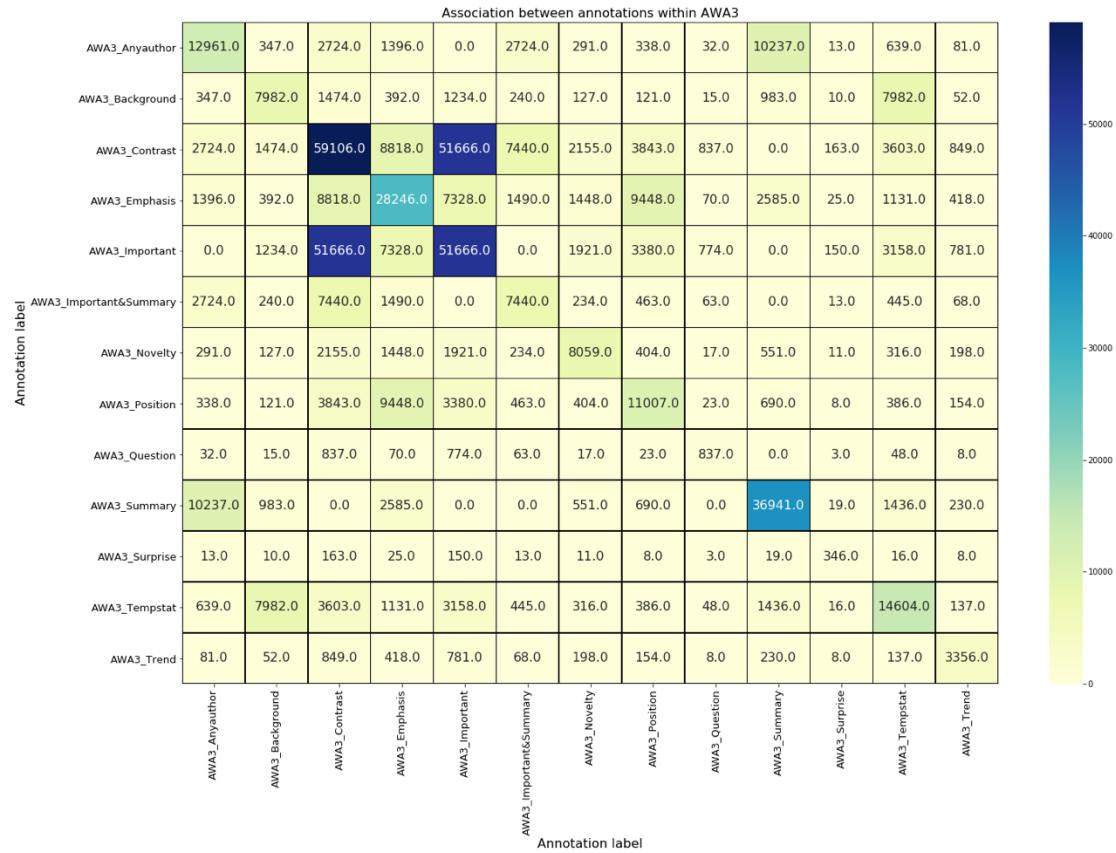
Since AWA tool annotated a sentence with different possible annotations, it is worth to see the association between these annotations. As shown in figure below, it seems that any sentence annotated with ‘Contrast’ will also be annotated with ‘Important’. Apart from that, most of the sentences annotated with ‘Emphasis’ will also be annotated with ‘Important’. Meanwhile, we also can see that ‘Position’ annotation does have strong association with not only ‘Important’ annotation but also ‘Emphasis’. Suppose two annotations are having very strong association, then either one is not necessary to be included as another separate annotation as it just cause redundancy.





## 9.0 Association within AWA3

In AWA3, there are two additional annotations called ‘AWA3\_Anyauthor’ and ‘AWA3\_Tempstat’ being considered. AWA3 also annotates each sentence with possible more than one annotations. The figures below show the association between the annotations within AWA3. One significant different from AWA is that the ‘Position annotation in AWA3 does have strong association with ‘Contrast’ annotation too.

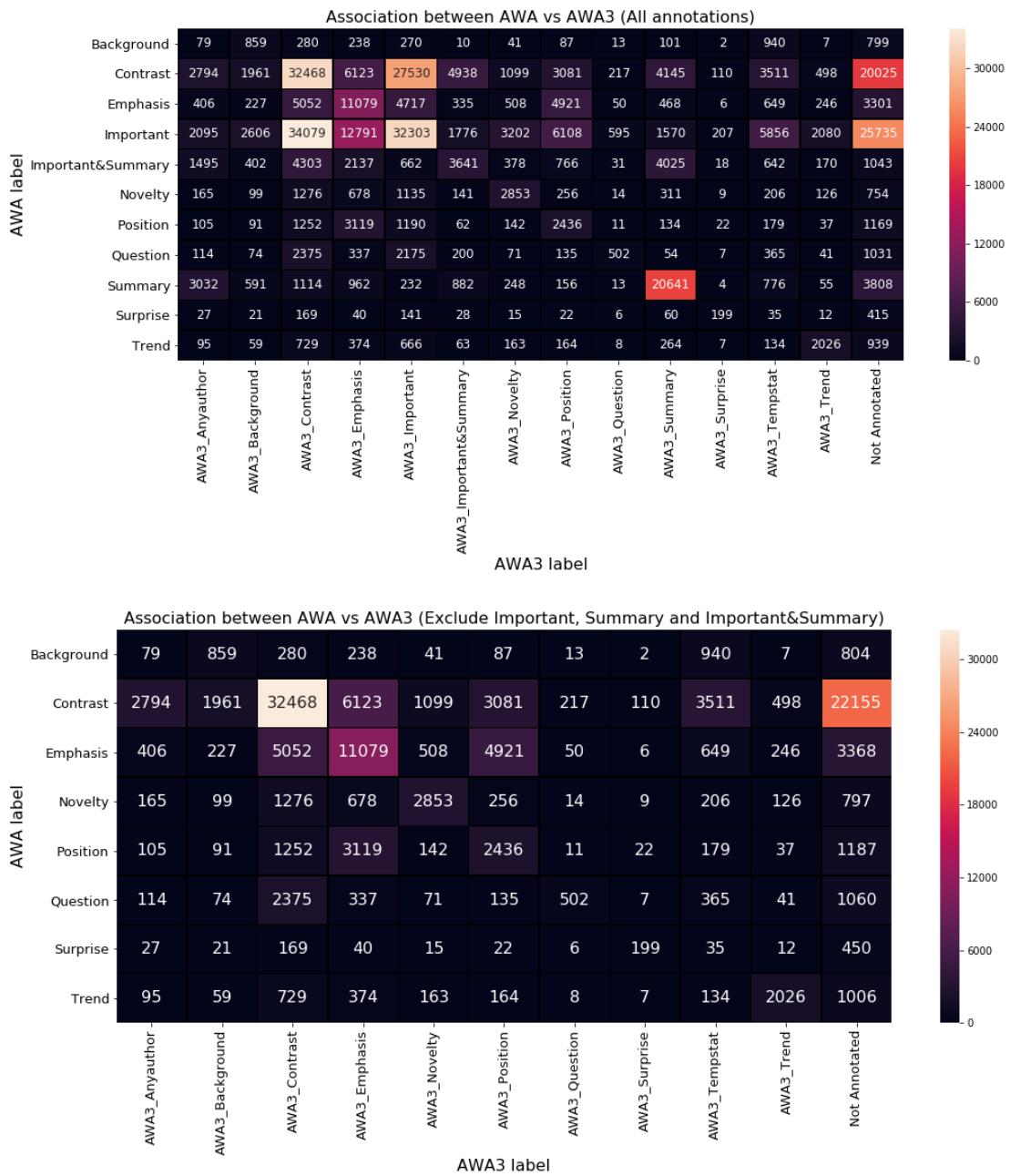


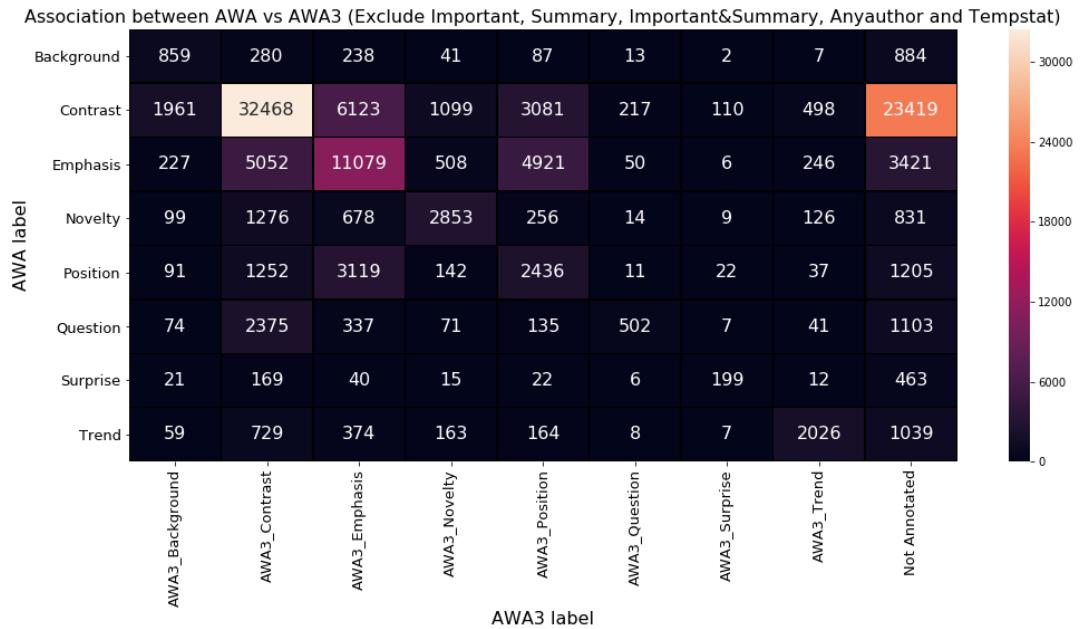
## 10.0 AWA and AWA3 Associations

AWA3 is somewhat supposed to be an improved version of AWA. So, it will be important to compare the association of these two tools.

### 10.1 AWA vs AWA3

Here, we will look at the association between AWA vs AWA3 in three different ways. One is association across all annotations, second is excluding the association ‘Important’, ‘Summary’ and ‘Important&Summary’ and lastly is the association excluding the ‘Important’, ‘Summary’, ‘Important&Summary’, ‘AWA3\_Anyauthor’ and ‘AWA3\_Tempstat’.

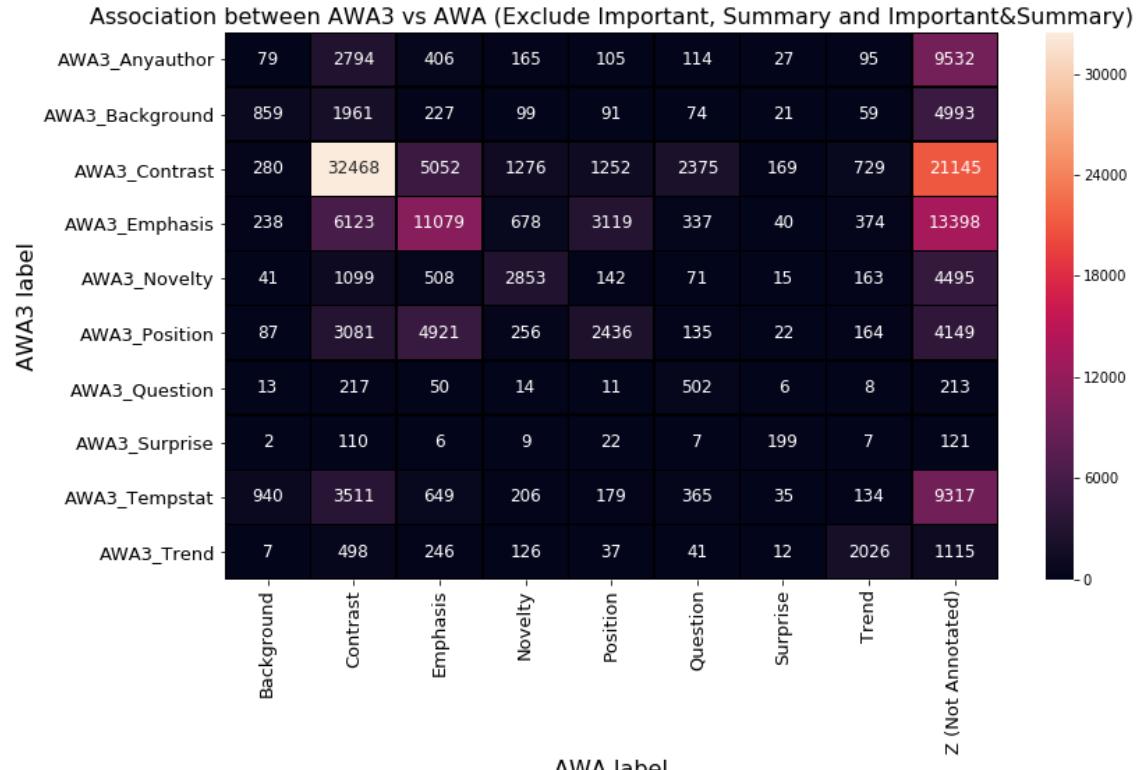
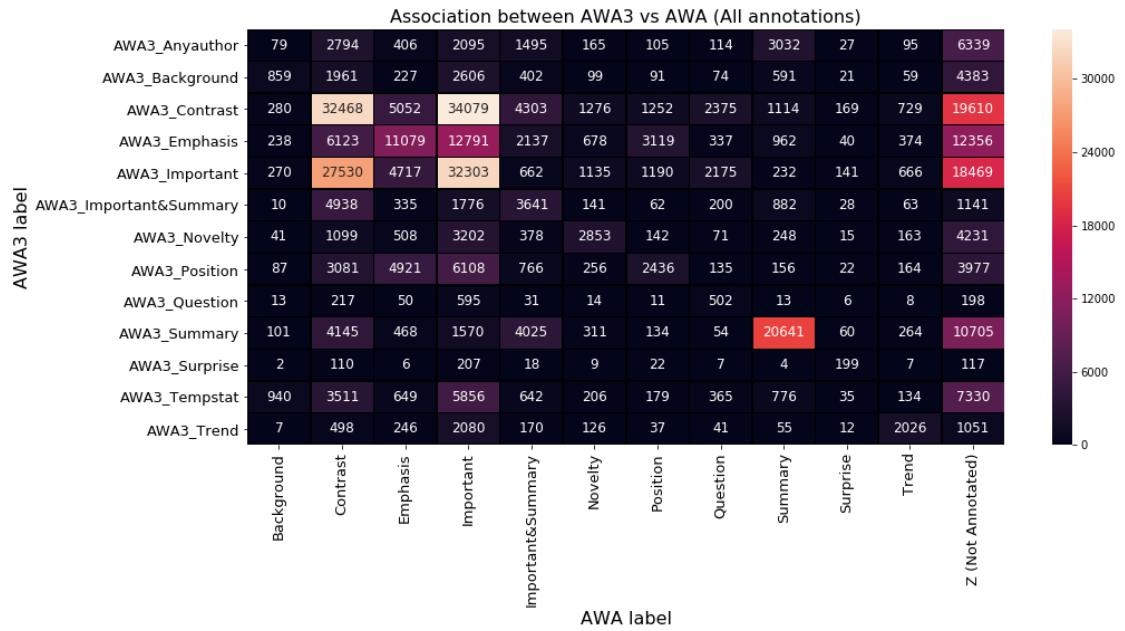


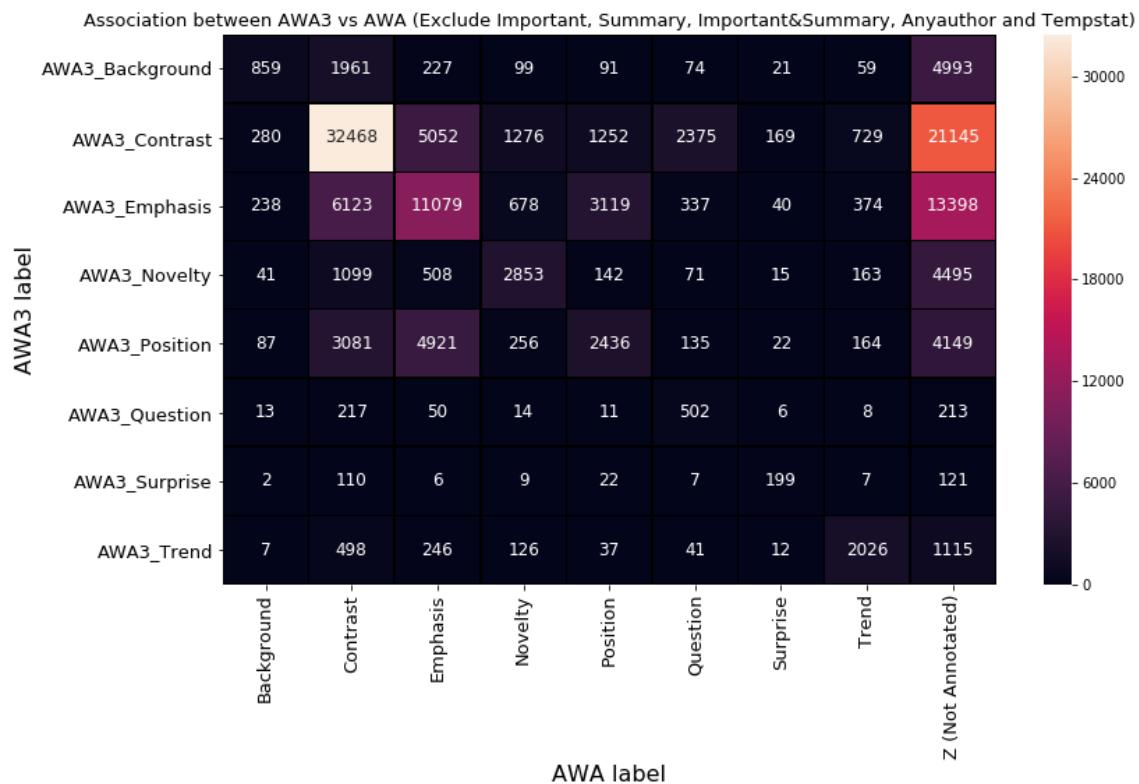


From the above diagrams, there are indeed some changes been made to AWA to become AWA3. As an example (diagram 3), there are many sentences (23419 sentences) that are annotated as ‘Contrast’ in AWA is no longer being annotated in AWA3. Apart from that, sentences annotated as ‘Question’ in AWA has no strong association with ‘AWA3\_Question’ but strong association with ‘AWA3\_Contrast’ as well as 1103 sentences in this annotation not being annotated. Apart from that, some sentences annotated with ‘Contrast’ in AWA have been shifted to be annotated as ‘Emphasis’ in AWA3. Meanwhile, some sentences annotated with ‘Emphasis in AWA have been shifted to be annotated as ‘Contrast’ in AWA3.

## 10.2 AWA3 vs AWA

Similar as 10.1 but this time is AWA3 vs AWA. In this way, we can see how many sentences that are annotated with AWA3 but not annotated with AWA.

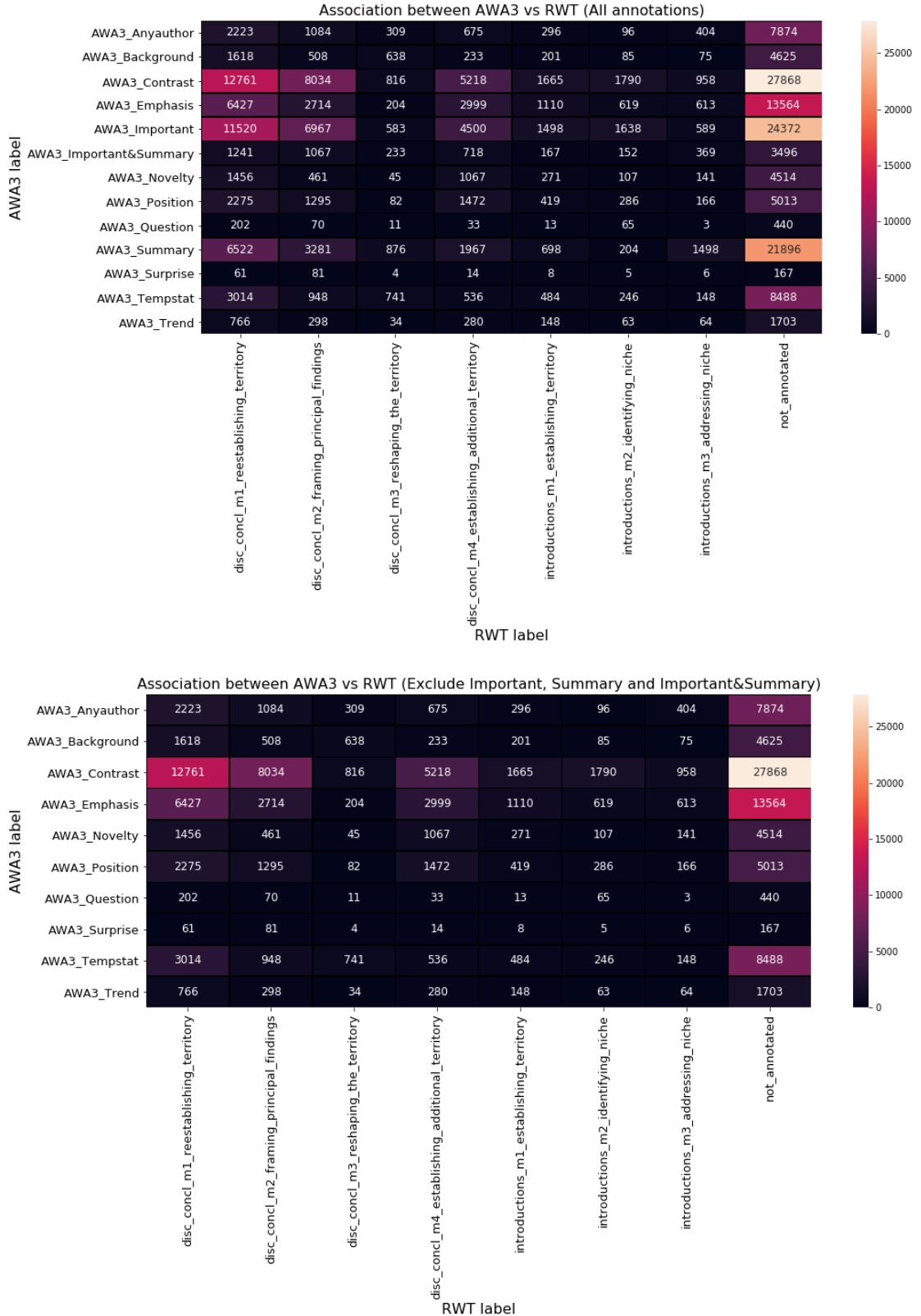


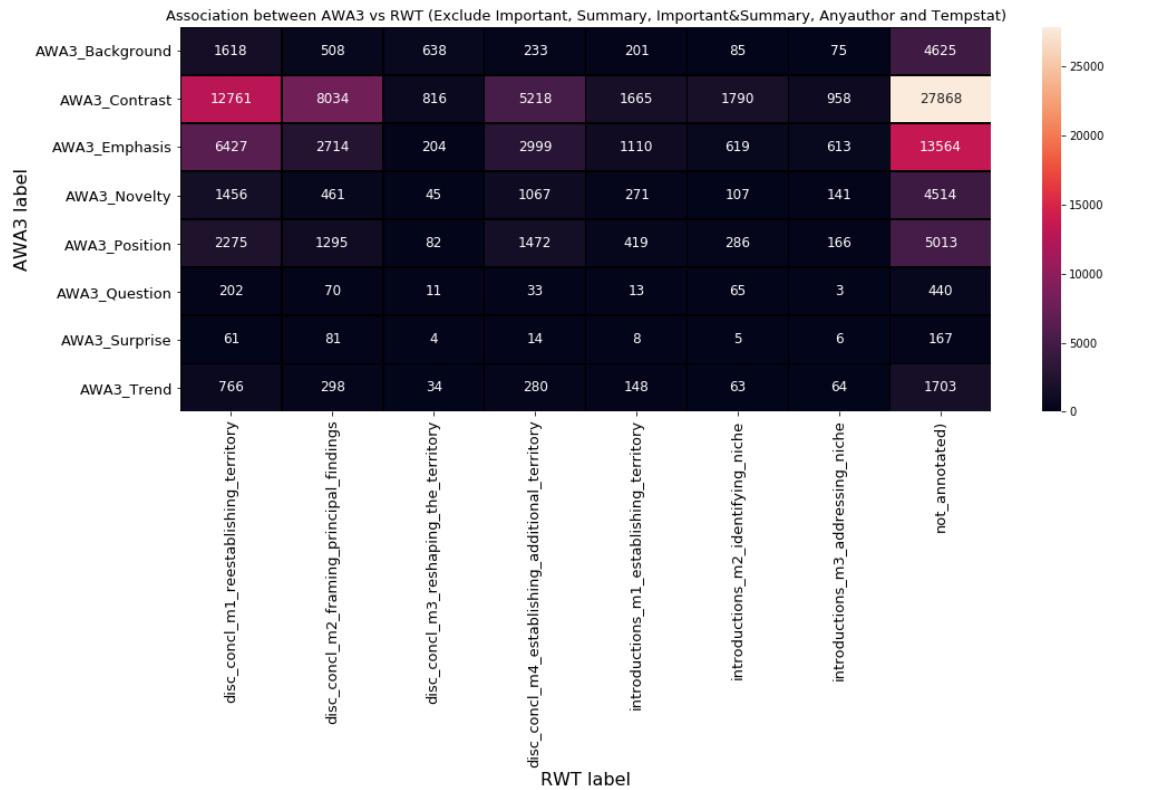


From the above diagrams, the obvious different between AWA3 vs AWA is that many sentences annotated by AWA3 are not annotated with any AWA annotation. That can be noticed in the last column of each of the diagrams.

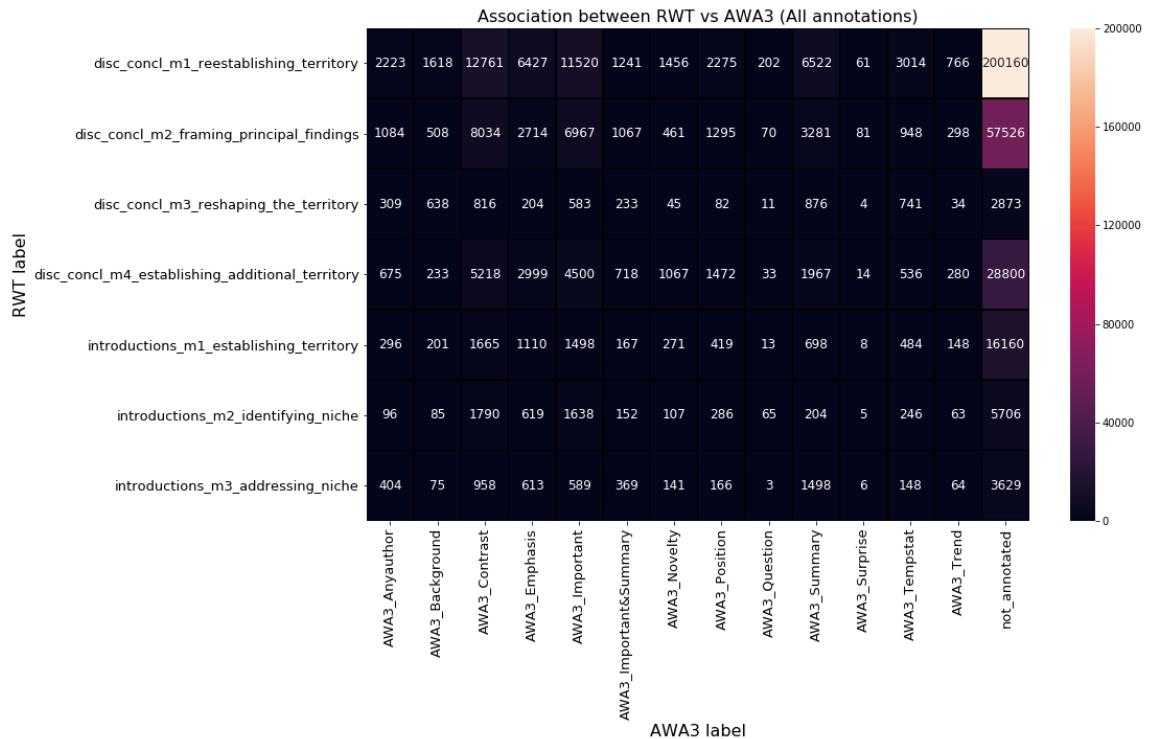
## 11. AWA3 and RWT associations

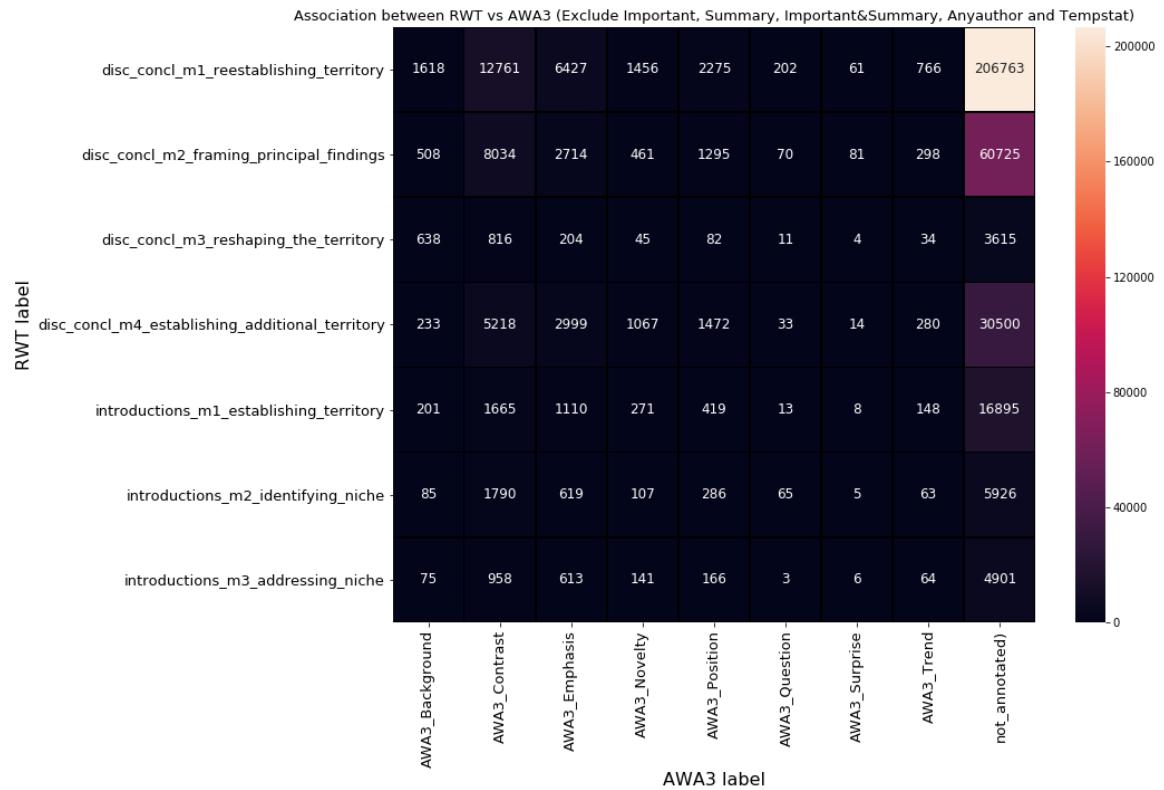
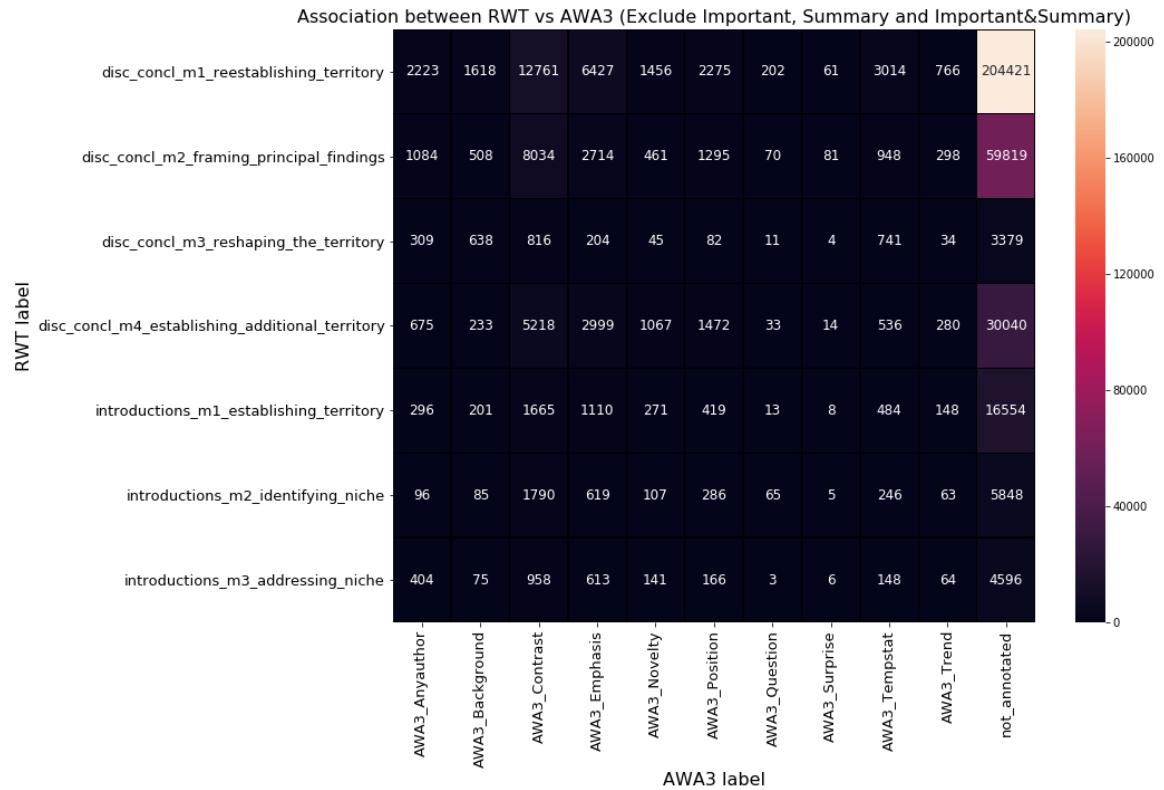
### 11.1 AWA3 VS RWT





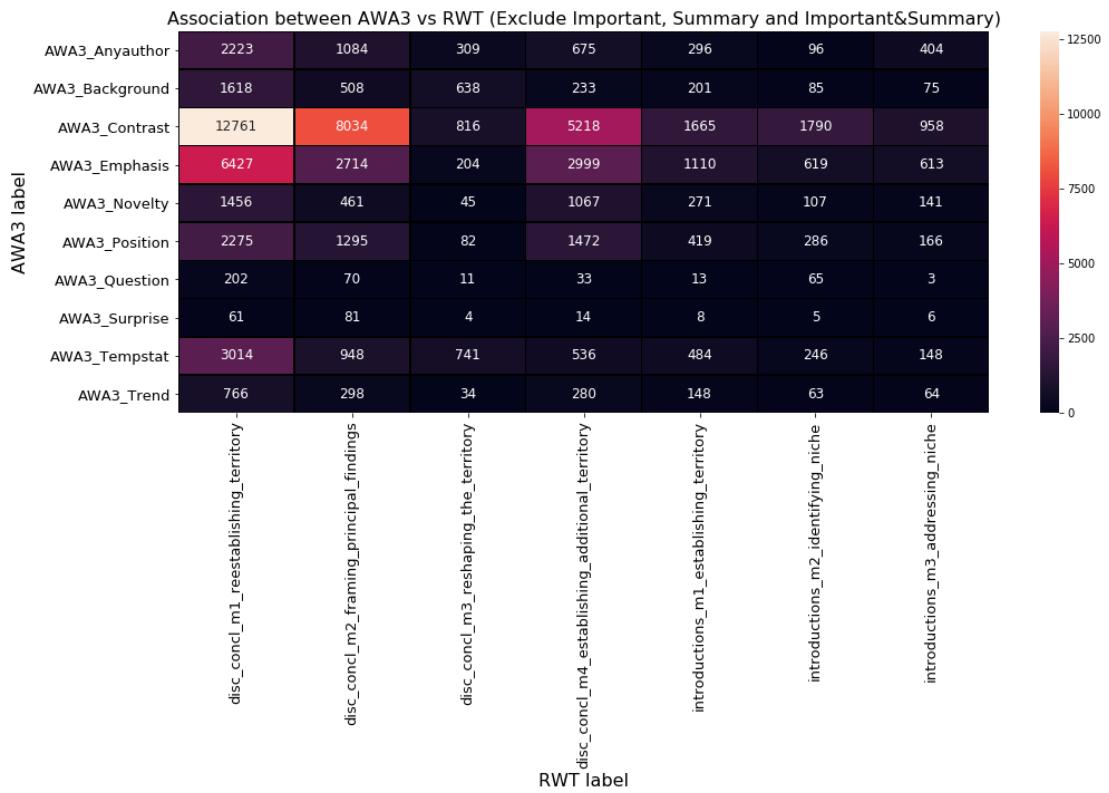
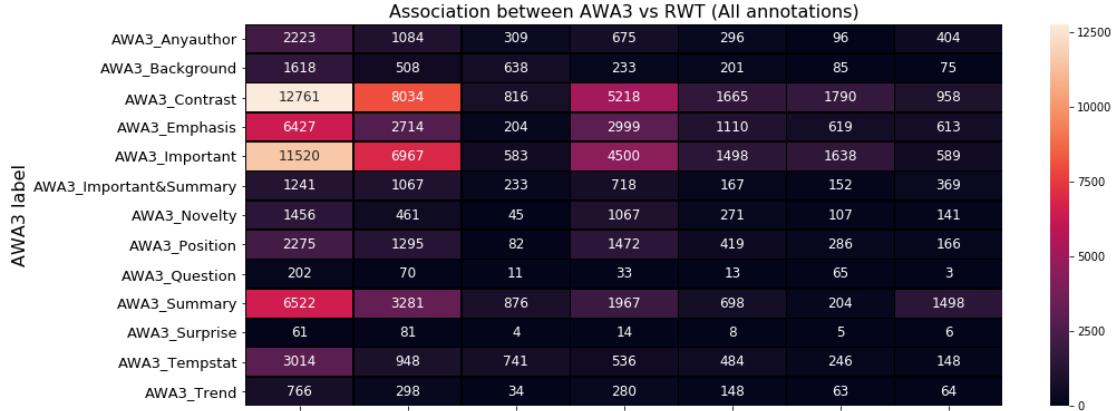
## 11.2 RWT VS AWA3





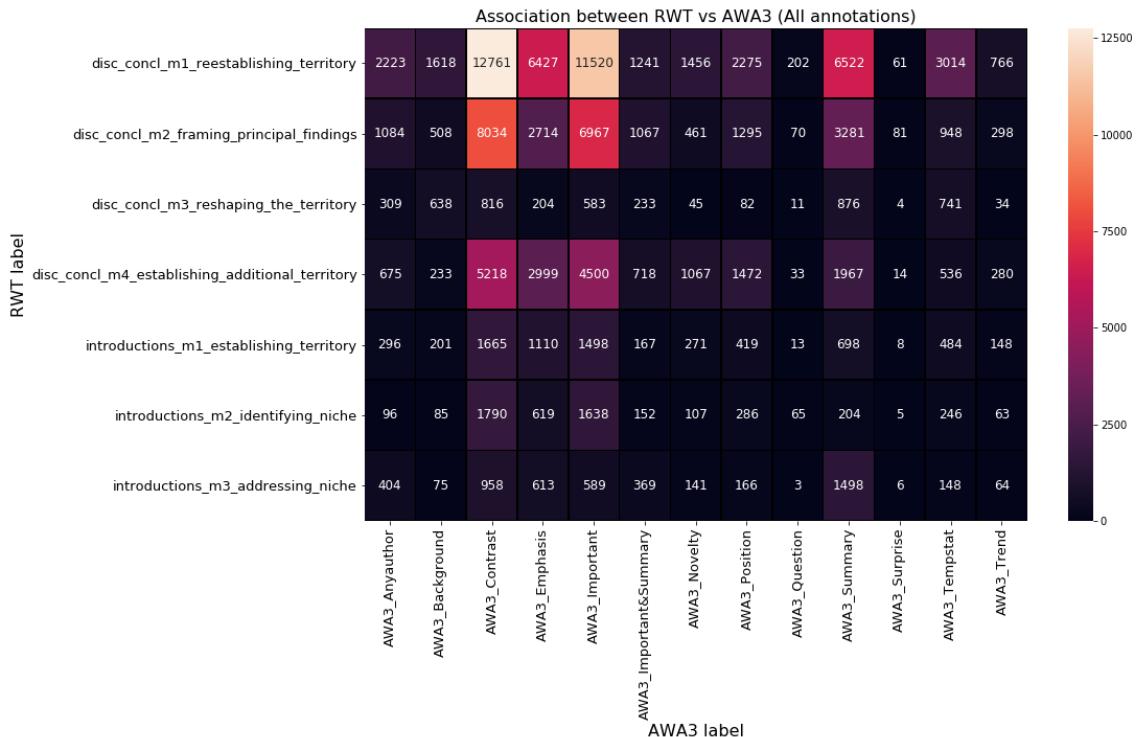
## 12. AWA3 and RWT associations (Ignoring sentences that not annotated)

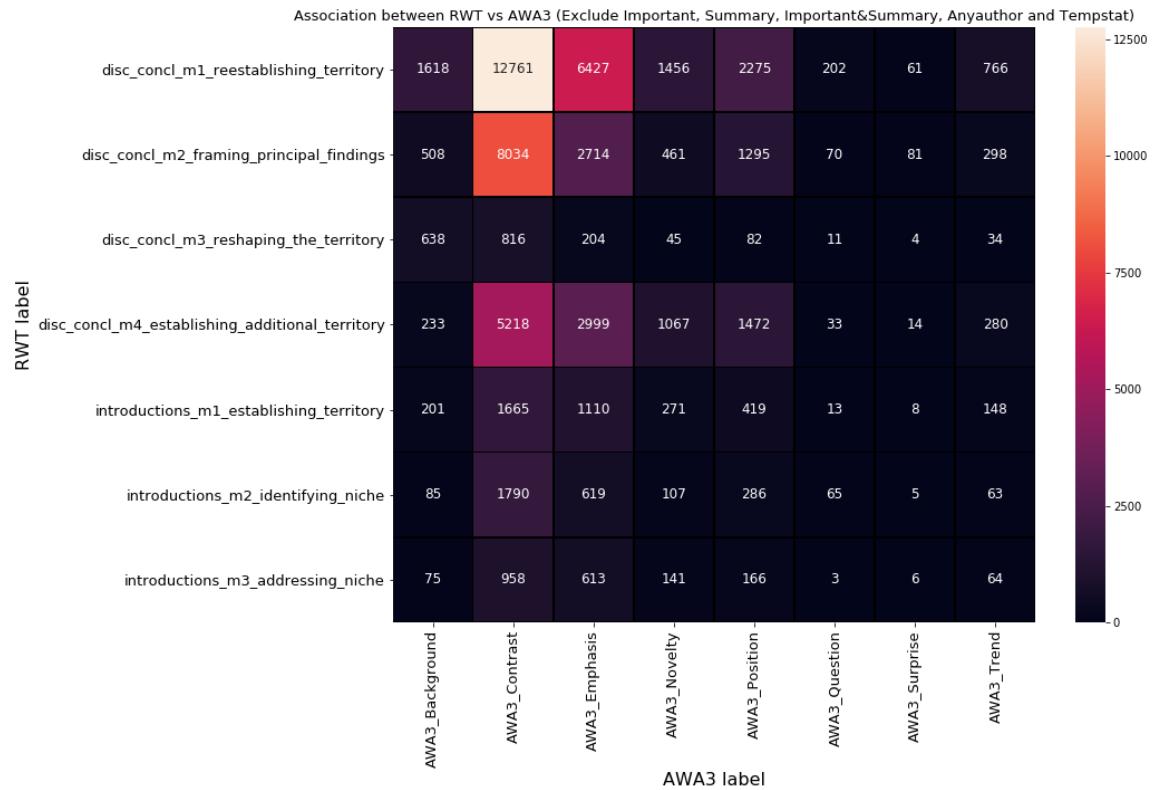
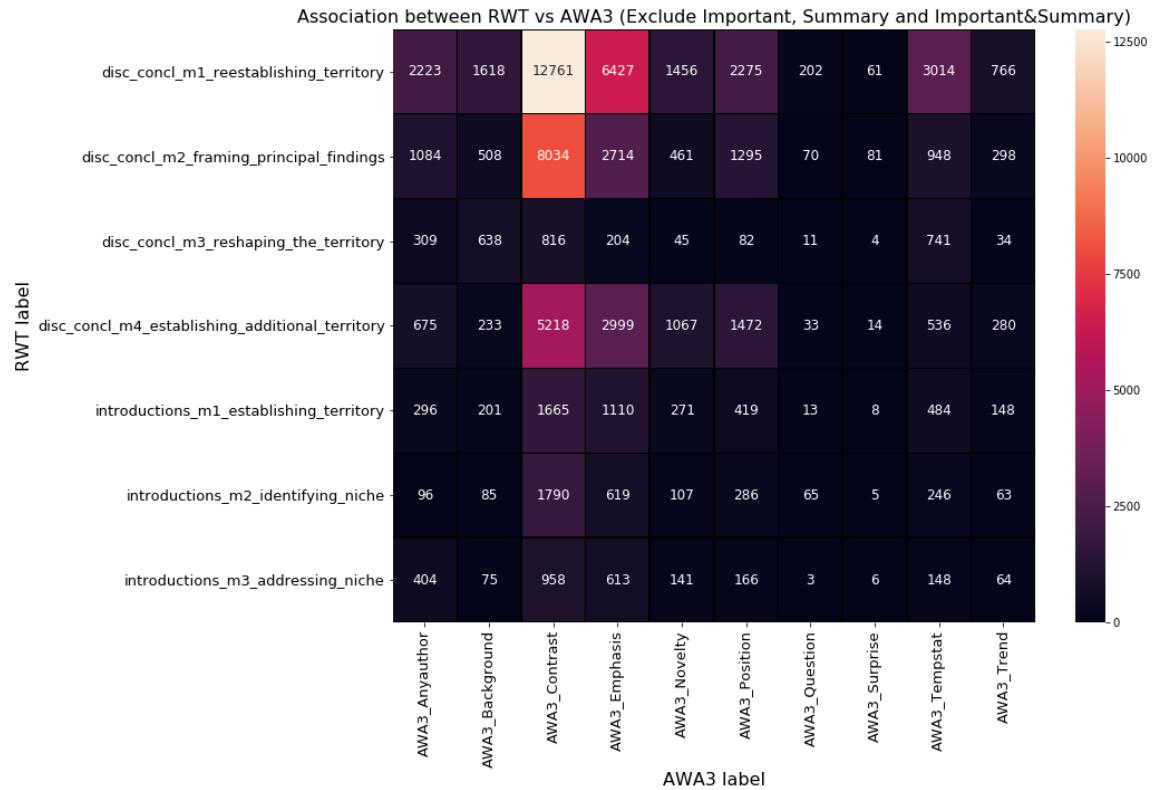
### 12.1 AWA3 VS RWT





## 12.2 RWT VS AWA3



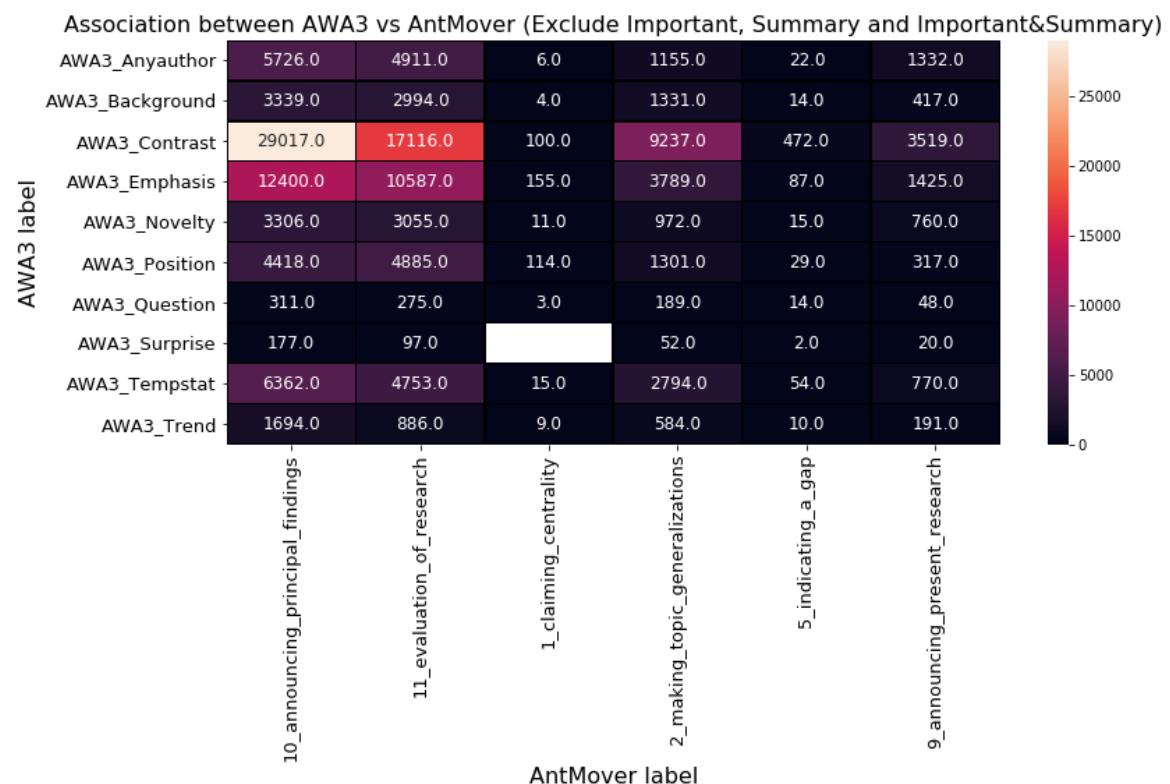
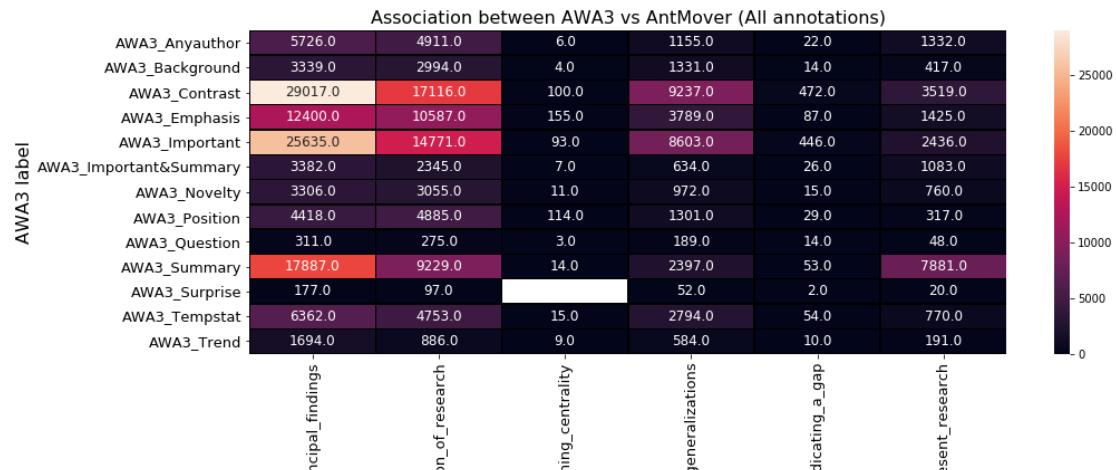


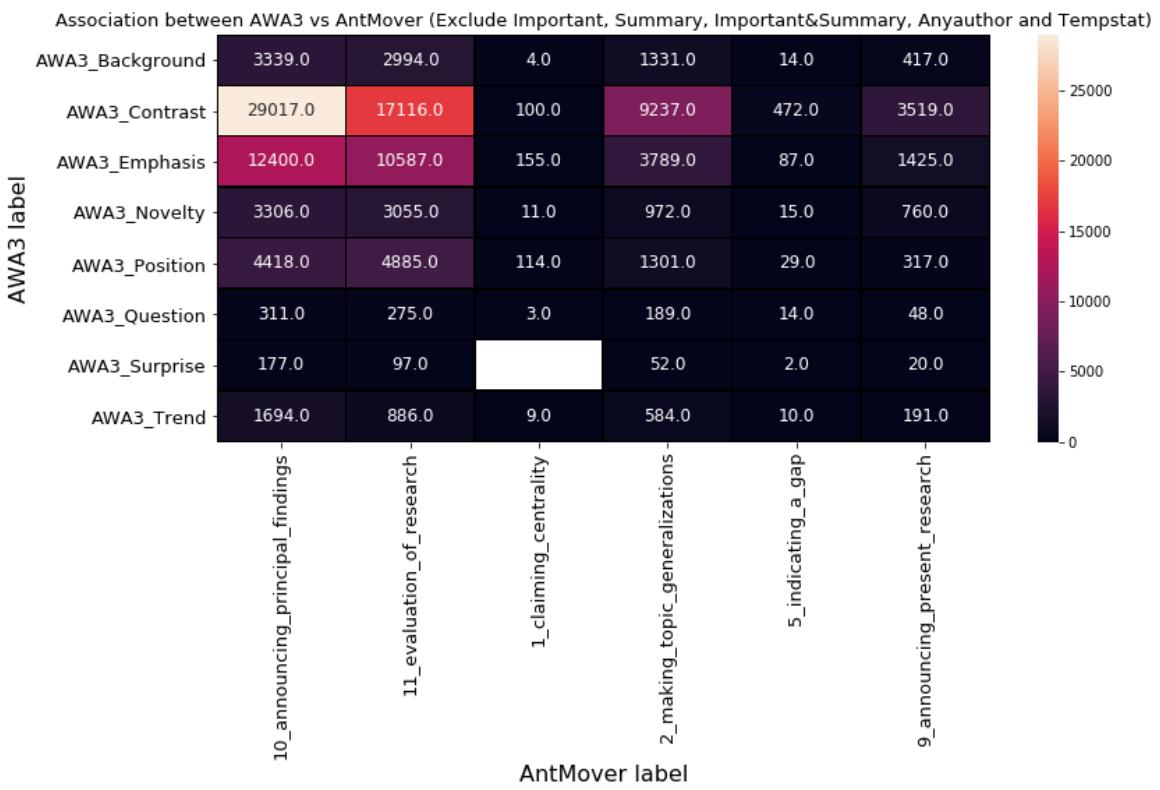
From the section 11 to 12, it is clearly to see that there is no direct association between AWA3 and RWT annotations. The only things can be derived from the above diagrams are some annotations such as ‘disc\_concl\_m1\_reestablishing\_territory’ and

'disc\_concl\_m2\_framing\_principle\_finding' are having high probability to be selected whenever that sentences are annotated as 'AWA3\_Contrast' and 'AWA3\_Empphasis' by AWA3.

## 13. AWA3 and AntMover associations

### 13.1 AWA3 VS AntMover

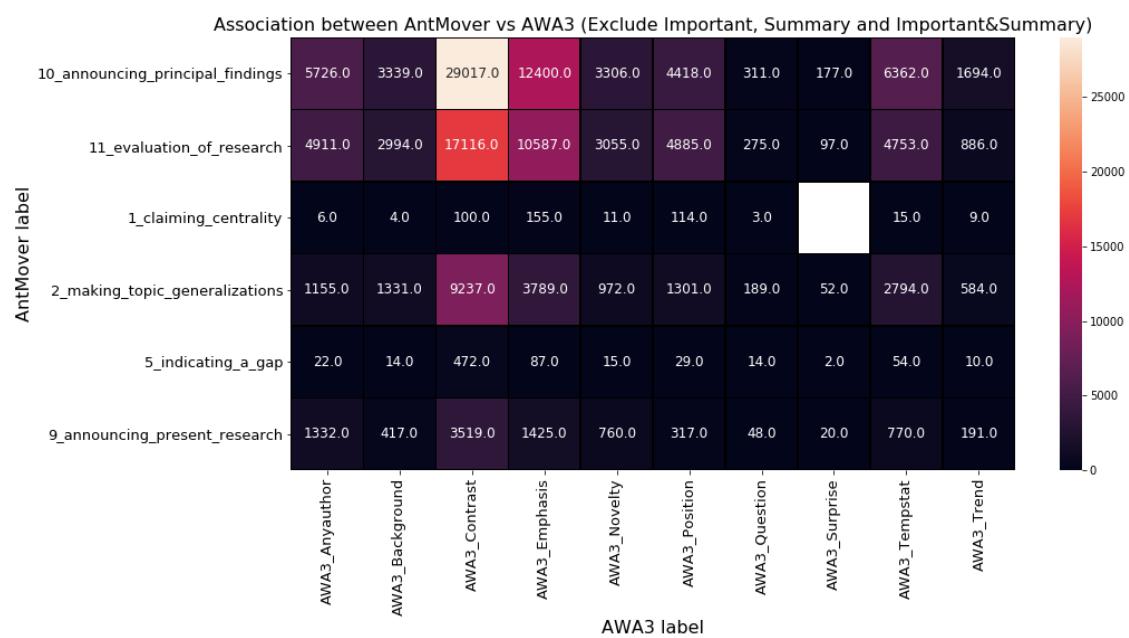
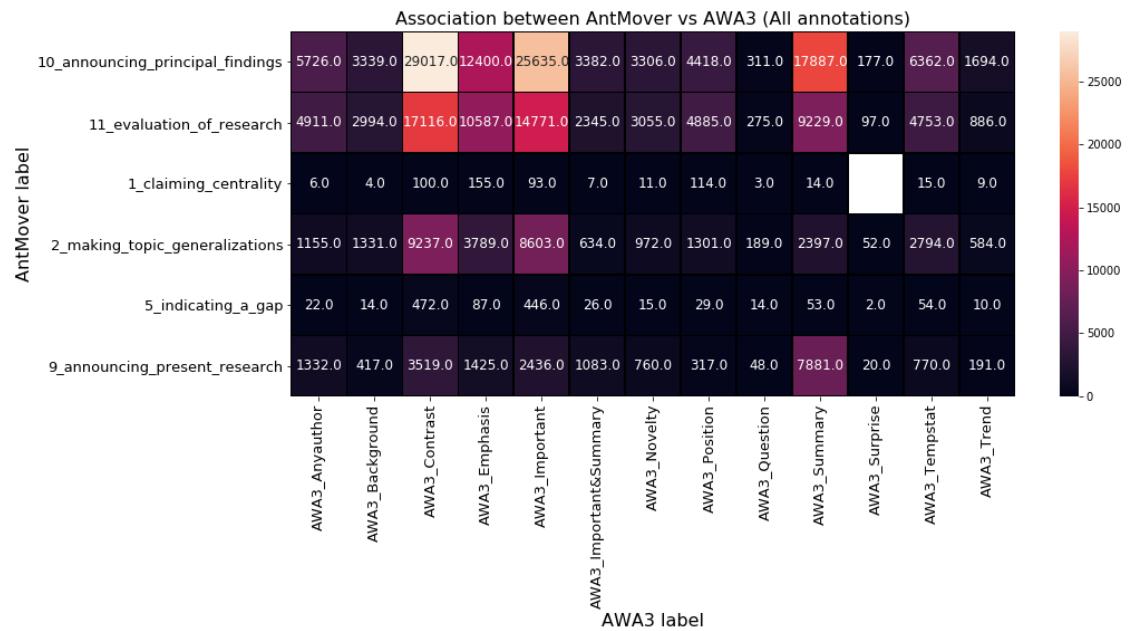


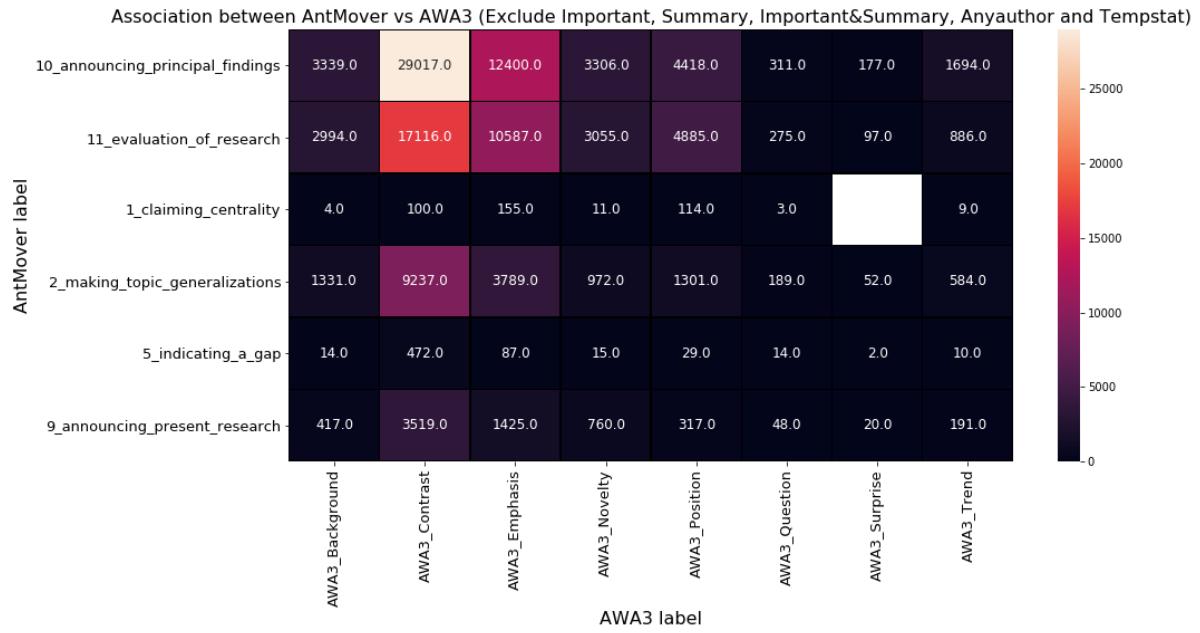


From the above three diagrams, we can see that there is no obvious association between AWA3 and AntMover annotations. However, there is one interesting thing to be noticed. There is no sentence being annotated as ‘AWA3\_Surprise’ is also being annotated as ‘1\_claiming\_centrality’.

### 13.2 AntMover VS AWA3

Since AntMover tool is annotating every sentence, we would not include sentences not annotated by AWA3 in this analysis or else it might mask some of the interesting facts from the analysis between the annotations of these tools.





As stated before, one interesting fact is that there is no sentence exist which is annotated as ‘1\_claiming\_centrality’ and ‘AWA3\_Surprise’.

## 14. Conclusion

Overall, there are some changes have been done on AWA3 from the previous version of AWA. However, these changes seem just a small changes as can see from the number of sentences being annotated as well as the comparison between AWA3 and AWA. Meanwhile, other tools such as RWT and AntMover seem do not have direct association to AWA3. There are some annotations from RWT or AntMover can be seen to have relations to certain annotations in AWA3 but these relations are not 1:1 relations. This report is not serve as a very detail statistical report but it is a report that provides the overall architecture of the database content as well as some basic comparison of the contents based on the tools. Therefore, this report can be significant for people who are intended to use the corpora in this database for research purpose.