

# Overview of Touché 2024: Argumentation Systems<sup>\*</sup>

Extended Version

Johannes Kiesel<sup>1,\*\*</sup>, Çağrı Çöltekin<sup>2</sup>, Maximilian Heinrich<sup>1</sup>, Maik Fröbe<sup>3</sup>, Milad Alshomary<sup>4</sup>, Bertrand De Longueville<sup>5</sup>, Tomaz Erjavec<sup>6</sup>, Nicolas Handke<sup>7</sup>, Matyáš Kopp<sup>8</sup>, Nikola Ljubešić<sup>6</sup>, Katja Meden<sup>6</sup>, Nailia Mirzhakhmedova<sup>1</sup>, Vaidas Morkevičius<sup>9</sup>, Theresa Reitis-Münstermann<sup>10</sup>, Mario Scharfbillig<sup>5</sup>, Nicolas Stefanovitch<sup>5</sup>, Henning Wachsmuth<sup>4</sup>, Martin Potthast<sup>11</sup> and Benno Stein<sup>1</sup>

<sup>1</sup>Bauhaus-Universität Weimar

<sup>2</sup>University of Tübingen

<sup>3</sup>Friedrich-Schiller-Universität Jena

<sup>4</sup>Leibniz University Hannover

<sup>5</sup>European Commission, Joint Research Centre (JRC)

<sup>6</sup>Jožef Stefan Institute

<sup>7</sup>Leipzig University

<sup>8</sup>Charles University

<sup>9</sup>Kaunas University of Technology

<sup>10</sup>Arcadia Sistemi Informativi Territoriali

<sup>11</sup>University of Kassel, hessian.AI, and ScaDS.AI

## Abstract

This paper is the extended overview of Touché: the fifth edition of the lab on argumentation systems that was held at CLEF 2024. With the goal to foster the development of support-technologies for decision-making and opinion-forming, we organized three shared tasks: (1) Human value detection (ValueEval), where participants detect (implicit) references to human values and their attainment in text; (2) Multilingual Ideology and Power Identification in Parliamentary Debates, where participants identify from a speech the political leaning of the speaker's party and whether it was governing at the time of the speech (new task); and (3) Image retrieval or generation in order to convey the premise of an argument with visually. In this paper, we describe these tasks, their setup, and participating approaches in detail.

## 1. Introduction

Decision-making and opinion-forming are everyday tasks, for which everybody has the chance to acquire knowledge on the Web on almost every topic. However, conventional search engines are

---

CLEF 2024: Conference and Labs of the Evaluation Forum, September 09–12, 2024, Grenoble, France

<sup>\*</sup>This overview extends the one published as part of the CLEF 2024 proceedings [1]

<sup>\*\*</sup>Corresponding author

✉ johannes.kiesel@uni-weimar.de (J. Kiesel); ccoltekin@sfs.uni-tuebingen.de (Ç. Çöltekin); maximilian.heinrich@uni-weimar.de (M. Heinrich); maik.froebe@uni-jena.de (M. Fröbe); m.alshomary@ai.uni-hannover.de (M. Alshomary); bertrand.de-longueville@ec.europa.eu (B. De Longueville); tomaz.erjavec@ijs.si (T. Erjavec); nicolas@bioinf.uni-leipzig.de (N. Handke); kopp@ufal.mff.cuni.cz (M. Kopp); nikola.ljubestic@ijs.si (N. Ljubešić); katja.meden@ijs.si (K. Meden); nailia.mirzhakhmedova@uni-weimar.de (N. Mirzhakhmedova); vaidas.morkevicius@ktu.lt (V. Morkevičius); theresa.reitis-munstermann@ext.ec.europa.eu (T. Reitis-Münstermann); mario.scharfbillig@ec.europa.eu (M. Scharfbillig); nicolas.stefanovitch@ec.europa.eu (N. Stefanovitch); h.wachsmuth@ai.uni-hannover.de (H. Wachsmuth); martin.pothast@uni-kassel.de (M. Potthast); benno.stein@uni-weimar.de (B. Stein)

🌐 <https://touche.webis.de> (Touché web page)

🆔 0000-0002-1617-6508 (J. Kiesel); 0000-0003-1031-6327 (Ç. Çöltekin); 0000-0001-5450-8203 (M. Heinrich); 0000-0002-1003-981X (M. Fröbe); 0000-0001-6142-9124 (M. Alshomary); 0009-0008-4215-5429 (B. De Longueville); 0000-0002-1560-4099 (T. Erjavec); 0000-0003-1349-4671 (N. Handke); 0000-0001-7953-8783 (M. Kopp); 0000-0001-7169-9152 (N. Ljubešić); 0000-0002-0464-9240 (K. Meden); 0000-0002-8143-1405 (N. Mirzhakhmedova); 0000-0002-2174-0396 (V. Morkevičius); 0009-0001-9828-1466 (T. Reitis-Münstermann); 0000-0002-1368-4187 (M. Scharfbillig); 0009-0000-2061-3216 (N. Stefanovitch); 0000-0003-2792-621X (H. Wachsmuth); 0000-0003-2451-0665 (M. Potthast); 0000-0001-9033-2217 (B. Stein)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

primarily optimized for returning *relevant* results, which is insufficient for collecting and weighing the pros and cons for a topic. To close this gap of technologies that support people in decision-making and opinion-forming, the Touché lab’s shared tasks<sup>1</sup> (<https://touche.webis.de>) call for the research community to develop respective approaches. In 2024, we organized the three following shared tasks:

1. Human Value Detection (a continuation of ValueEval’23 @ SemEval [2]) features two subtasks in ethical argumentation of detecting human values in texts and their attainment, respectively.
2. Ideology and Power Identification in Parliamentary Debates features two subtasks in debate analysis of detecting the ideology and position of power of the speaker’s party, respectively (new task).
3. Image Retrieval/Generation for Arguments (third edition, now joint task with ImageCLEF) is about the retrieval or generation of images to help convey an argument’s premise.

In total, 20 teams participated in Touché in 2024. Nine teams participated in the human value detection task (cf. Section 4)—of which six submitted a notebook paper—and submitted 21 runs. Most teams integrated DeBERTa [3], RoBERTa [4], or the multi-lingual XLM-RoBERTa [5]. Only one team employed a generative approach (employing GPT-4o). Nine teams participated in the multilingual ideology and power identification task (cf. Section 5) and submitted 52 runs. The majority of teams participated in both subtasks. While traditional machine learning methods like support vector classifiers or logistic regression with n-gram features were more common among participating teams, higher-scores were typically obtained by teams using pretrained models. Two teams participated in the image retrieval/generation for arguments task (cf. Section 6) and submitted 8 runs. Both teams used similarity embeddings between images and text. One team used CLIP [6], the other a DPR [7] inspired approach. The corpora, topics, and judgments created at Touché are freely available to the research community on the lab’s website.<sup>2</sup>

## 2. Related Work

Argumentation systems are diverse and are connected to many fields within and outside of computer science. The following sections review the related work for each Touché task of 2024.

### 2.1. Human Value Detection

Due to their outlined importance, human values have been studied both in the social sciences [8] and in formal argumentation [9] for decades. According to the former, a “value is a (1) belief (2) pertaining to desirable end states or modes of conduct, that (3) transcends specific situations, (4) guides selection or evaluation of behavior, people, and events, and (5) is ordered by importance relative to other values to form a system of value priorities.” For cross-cultural analysis, Schwartz derived 48 value questions from universal individual and societal needs, including concepts such as *obeying all the laws* and *being humble* [10]. Based on these taxonomies are several studies in the social sciences, which could greatly benefit from the automated methods our task aims at [11]. See Scharfbillig et al. [12] for a recent overview and practical insights from the social sciences.

Moreover, several works in computer science utilize values. For example, in the context of interactive systems, to tune interactive chat-based agents or texts in general towards morally acceptable behavior [13, 14]. A related dataset is ValueNet [15], which contains 21K one-sentence descriptions of social scenarios (taken from SOCIAL-CHEM-101 [16]) annotated for the 10 value categories of an earlier version of Schwartz’ value taxonomy. A major difference to the Touché24-ValueEval dataset are the more ordinary situations in ValueNet (e.g., whether to say “I miss mom”). Our earlier work analyzed values in short arguments [17, 2].

---

<sup>1</sup>‘touché’ confirms “a hit in fencing or the success or appropriateness of an argument, an accusation, or a witty point.” [<https://merriam-webster.com/dictionary/touche>]

<sup>2</sup><https://touche.webis.de/>

## 2.2. Ideology and Power Identification

Parliamentary data has a high societal impact and provides publicly available sources for analyzing (argumentative) language. Thus the number of resources based on parliamentary proceedings [18, 19], and computational and linguistics analyses of parliamentary debates [20, 21] increased in recent years.

The present task is about two important aspects of the political discourse, *ideology* and *power*. Although a simplification, political orientation on the left-to-right spectrum has been one of the defining properties of political ideology [22, 23]. Power is another factor that shapes the political discourse [24, 25, 26]. Automatic identification of political orientation from texts has attracted considerable interest [27, 28, 29, 30, 31], including a few recent shared tasks [32, 33]. The present task differs from the earlier ones, with respect to the source material (parliamentary debates, rather than the popular sources of social media or news) and multilinguality. Despite its central role in critical discourse analysis, to the best of our knowledge, power in parliamentary debates has not been studied computationally. There has been only a few recent computational studies providing indications of linguistic differences between governing and opposition parties [34, 35, 36, 37]. The present shared task and associated data is likely to provide a reference for the future studies investigating power in political discourse.

## 2.3. Image Retrieval/Generation for Arguments

Images are a powerful tool for visual communication. They can provide contextual information and express, underline, or popularize an opinion [38], thereby taking the form of subjective statements [39]. Some images express both a premise and a conclusion, making them full arguments [40, 41]. Other images may provide contextual information only and have to be combined with a textual conclusion to form a complete argument. In this regard, a recent SemEval task distinguished a total of 22 persuasion techniques in memes alone [42]. Moreover, argument quality dimensions like acceptability, credibility, emotional appeal, and sufficiency [43] all apply to arguments that include images as well.

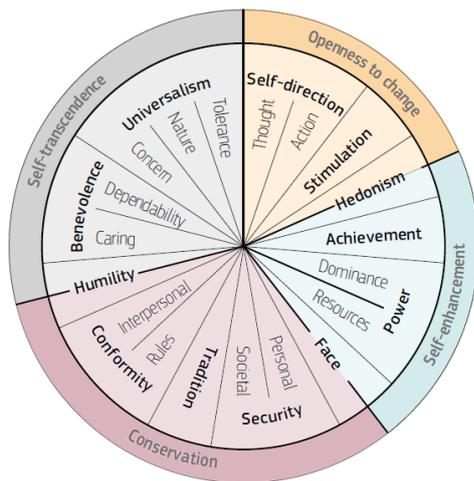
## 3. Lab Overview and Statistics

For the fifth edition of the Touché lab, we received 68 registrations from 22 countries (vs. 41 registrations in 2023). The most lab registrations came from India (24). Out of the 68 registered teams, 20 actively participated in this year’s Touché edition (9, 9, and 2 teams submitting valid runs for Task 1, 2, and 3, respectively). Active teams in previous editions were: 7 in 2023, 23 in 2022, 27 in 2021, and 17 in 2020.

We used TIRA [44] as the submission platform for Touché 2024 through which participants could either submit code, software, or run files.<sup>3</sup> Code and software submissions increase reproducibility, as the software can later be executed on different data of the same format. To submit software, a team implemented their approach in a Docker image that they then uploaded to their dedicated Docker registry in TIRA. Software submissions in TIRA are immutable, and after the docker image had been submitted, the teams specified the to-be-executed command—the same Docker image can thus be used for multiple software submissions (e.g., by changing some parameters). A team could upload as many Docker images or software submissions as they liked; only they and TIRA had access to their dedicated Docker image registry (i.e., the images were not public while the shared task was ongoing). To improve reproducibility, TIRA executes software in a sandbox by removing the internet connection (ensuring that the software is fully installed in the Docker image which eases rerunning software later, as libraries and models must be installed in an image). For the execution, participants could select the resources that their software had available for execution, from 1 CPU core with 10 GB RAM up to 5 CPU cores with 50 GB RAM and 1 Nvidia A100 GPU with 40 GB RAM. Participants could run their software multiple times using different resources to study the scalability and reproducibility (e.g., whether the software executed on a GPU yields the same results as on a CPU). TIRA used a Kubernetes cluster with 1,620 CPU cores, 25.4 TB RAM, 24 GeForce GTX 1080 GPUs, and 4 A100 GPUs to schedule and execute the software submissions, to allocate the resources that the participants selected.

---

<sup>3</sup><https://tira.io>



Inner circle: 19 human values  
(see <https://valueeval.webis.de>)

Outer circle: four motivational directions  
(not used in this task)

- **Openness to change**  
Being independent and exploring
- **Self-enhancement**  
Seeking pleasure, wealth, and esteem
- **Conservation**  
Preserving group cohesion, order, and security
- **Self-transcendence**  
Helping others, close ones, and nature

**Figure 1:** The 19 values used in this task, shown in the Schwartz value taxonomy [10].

## 4. Task 1: Human Value Detection (ValueEval'24)

The goal of this task is to develop approaches that allow for the large-scale analysis of human values behind texts. In argumentation, one has to consider that people have different beliefs and priorities of what is generally worth striving for (e.g., personal achievements vs. humility) and how to do so (e.g., being self-directed vs. respecting traditions), referred to as (human) values. By analyzing corpora of texts, for example for news portals or political parties, one can develop an understanding of the values that the authors deem the most important.

### 4.1. Task Definition

The task is to identify the values of the widely accepted value taxonomy of Schwartz [10] (cf. Figure 1) and their attainment in long texts of nine languages (Bulgarian, Dutch, English, French, German, Greek, Hebrew, Italian, and Turkish). This taxonomy has been replicated in over 200 samples in 80 countries and is the backbone of value research [12]. A value can either be mentioned as something that is or should be attained (i.e., lead towards fulfilling the value) or something that is constrained, i.e., not attained. For example, for Security, (partial) attainment would mean that something is made safer or healthier. In contrast, an event can be stated in a way that thwarts or constrains safety or health. Participating teams can submit software in one or both of two subtasks: (1) Given a text, for each sentence, detect which human values the sentence refers to; and (2) Given a text, for each sentence and value this sentence refers to, detect whether this reference (partially) attains or constrains the value.

### 4.2. Data Description

The task employs a collection of 2648 human-annotated texts in nine languages from news articles and political manifestos. Texts are sampled to reflect diverse opinions (different parties; mainstream news and others) from 2019 to 2023. The data is annotated as part of the ValuesML project<sup>4</sup> by over 70 value scholars. The annotators marked segments in the texts, selected from 19 values the values that the segment refers to most, and selected for each of these values whether the segment (partially) attains or constrains the value, or whether attainment is unclear. Dedicated team leaders per language trained the respective annotators, discussed sentences for which annotators disagreed in their teams, and consolidated annotations into one ground truth. The team leaders discussed issues with us in bi-weekly meetings. Moreover, we discussed with the team leaders the current holistic inter-annotator agreement [45] and its change compared to the previous meeting to monitor annotation

<sup>4</sup>[https://knowledge4policy.ec.europa.eu/projects-activities/valuesml-unravelling-expressed-values-media-informed-policy-making\\_en](https://knowledge4policy.ec.europa.eu/projects-activities/valuesml-unravelling-expressed-values-media-informed-policy-making_en)

**Table 1**

Overview of the Touché24-ValueEval dataset by language, with the respective number of texts, sentences, annotator agreement as measured by Krippendorff’s  $\alpha$ , and the thousandths of these sentences with any or a specific value (attained or constrained). Languages are Bulgarian (BE), German (DE), Greek (EL), English (EN), French (FR), Hebrew (HE), Italian (IT), Dutch (NL), and Turkish (TR).

Lang.	Texts	Sentences	$\alpha$	Sentences with value (%)																			
				Any value	Self-direction: thought	Self-direction: action	Stimulation	Hedonism	Achievement	Power: dominance	Power: resources	Face	Security: personal	Security: societal	Tradition	Conformity: rules	Conformity: interpersonal	Humility	Benevolence: caring	Benevolence: dependability	Universalism: concern	Universalism: nature	Universalism: tolerance
BG	260	6 919	0.495	641	010	055	046	005	075	053	053	021	011	108	020	089	009	002	059	021	071	023	005
DE	261	9 183	0.367	533	018	055	034	011	079	032	038	020	026	059	009	072	015	002	017	015	050	026	014
EL	328	7 349	0.696	615	003	013	029	003	054	074	089	018	011	130	006	060	046	000	024	032	054	025	014
EN	408	10 305	0.409	306	004	025	005	004	043	016	016	006	014	053	008	036	016	003	006	007	031	012	008
FR	219	4 650	0.685	304	005	023	016	005	019	024	015	020	021	065	006	030	010	001	012	007	038	020	009
HE	250	7 331	0.557	859	025	042	021	003	081	122	094	032	029	170	031	096	011	002	016	041	080	022	015
IT	276	6 379	0.610	632	010	015	072	008	133	053	082	029	013	071	003	076	002	000	018	004	045	038	009
NL	323	10 982	0.411	366	014	029	004	003	039	030	037	010	009	072	004	033	005	002	004	017	043	019	009
TR	323	11 133	0.463	473	015	046	027	022	059	025	045	016	042	072	027	071	007	004	047	025	036	014	007
All	2 648	74 231	0.546	512	012	035	026	008	063	045	050	018	020	086	013	061	013	002	022	019	048	021	010

quality and coherence across documents and languages. To measure annotator agreement, we computed Krippendorff’s  $\alpha$  before curation for all language teams individually and overall (cf. Table 1). We see this agreement as sufficient, and the curation process increased the annotation quality even further.

For Touché, the dataset is automatically split into sentences using Trankit version 1.1.1 [46]. Table 2 shows the dataset format. The dataset is provided both in the original language and automatically translated to English, either using DeepL or, for Hebrew, Google Translate.<sup>5</sup> The dataset is split into sets by texts, so that 60% / 20% / 20% of sentences are in the training / validation / test set, respectively.<sup>6</sup>

Table 1 shows the size and value distribution for each language. The number of texts per language are between 219 (French) and 408 (English). The number of sentences per language are between 4 650 (French) and 11 133 (Turkish). Only 30.4% of the French sentences are annotated as referring to a value, but 85.9% of Hebrew sentences. The value frequency is between 0.2% (Humility) and 8.6% (Security: societal). This in-balance between languages and values makes the problem especially challenging.

### 4.3. Participant Approaches

In 2024, nine teams participated in this task (of which six submitted a notebook paper) and submitted 21 runs. Moreover, we added two baseline runs for comparison. Five of the six teams that submitted a paper relied on DeBERTa [3], RoBERTa [4], or the multi-lingual XLM-RoBERTa [5]. The other team (Eric Fromm) used GPT-4o.<sup>7</sup> Two teams work with the multi-lingual dataset (Arthur Schopenhauer, Hierocles of Alexandria) whereas the others use the English translations only. Only one team (Hierocles of Alexandria) used the sentence sequence, whereas the other teams classified each sentence individually.

<sup>5</sup><https://www.deepl.com/pro-api> and <https://cloud.google.com/translate>

<sup>6</sup>Dataset: <https://zenodo.org/doi/10.5281/zenodo.10396293>

<sup>7</sup><https://openai.com/index/hello-gpt-4o/>

**Table 2**

Excerpt of the dataset for the human value detection task. The dataset comes in six directories: training, validation, and test data for both the original multi-lingual dataset and its automatic translation to English. Each directory contains a `sentences.tsv` where each row corresponds to one sentence. The training and validation directories also each contain a `labels.tsv` where each row corresponds to a sentence in `sentences.tsv` and columns 3–40 correspond to labels (attained and constrained for each of the 19 values). Label values in the `labels.tsv` are either 1.0 if the sentence refers to that value and attainment polarity, 0.0 if it does not, or 0.5 if the sentence refers to that value but the attainment polarity is unclear (0.2% of cases).

`sentences.tsv` (3 columns)

Text-ID	Sentence-ID	Text
EN_012	1	Who designed global guidelines for puberty blockers?
EN_012	2	More and more children and young people believe they have to question their ...
EN_012	3	Some 60 minors were treated in the Netherlands in 2010, but has increased to ...

`labels.tsv` (40 columns)

Text-ID	Sentence-ID	Self-direction: thought attained	Self-direction: thought constrained	...
EN_012	1	0.0	0.0	...
EN_012	2	1.0	0.0	...
EN_012	3	0.0	0.0	...

**Baselines.** We provide two baselines, that also served to kickstart the participants’ approaches:<sup>8</sup> (1) a random baseline assigns per sentence a uniformly random value “confidence” to each value in subtask 1 and randomly distributes this confidence between attained and constrained for subtask 2; and (2) a BERT [47] baseline trained for multi-label classification for all 38 combinations of value and attainment.

**Team Arthur Schopenhauer [48].<sup>9</sup>** The team used the multi-lingual dataset and analyzed the sentences independently. They approached subtask 1 as a classification problem. A *no-label* class was added for sentences without assigned value, and sentences with *Humility* were ignored due to the scarcity of that value. The 6% of sentences with more than one assigned value were ignored, as well. Different models were fine-tuned for English texts (deberta-v2-xxlarge [3]) and others (xlm-roberta-large [5]). In both cases, an ensemble with a thresholded soft voting scheme of four models was employed: one model for each combination of two seeds and two loss functions. For loss functions the authors report that cross entropy lead to higher results in their preliminary tests for frequent values but weighted cross entropy did so for infrequent values. The team approached subtask 2 as a binary classification problem, ignoring the few sentences with *unknown* attainment. Their approach is otherwise the same as for subtask 1, except that only a single model was employed instead of an ensemble (with cross entropy loss) based on results from their preliminary tests.

**Team Edward Said [49].** The team used the English translations of the dataset and analyzed the sentences independently. To counter the label imbalance, the team upsampled sentences by a factor of four if the associated label is one of 14 underrepresented labels (value + attainment; out of 38). They selected the 14 labels that are infrequent in total or in comparison to the label for the same value with other attainment. They fine-tuned a RoBERTa [4] and DeBERTa [3] model for multi-label classification.

**Team Eric Fromm [50].** The team used the English translations of the dataset and analyzed the sentences independently. They employed GPT-4o for zero-shot classification, prompting with the annotator guide’s 19 value descriptions to select at most one per sentence. They did not tackle subtask 2.

<sup>8</sup><https://github.com/touche-webis-de/touche-code/tree/main/clef24/human-value-detection/approaches>

<sup>9</sup>Code: <https://github.com/h-uns/clef2024-human-value-detection>

Models: <https://huggingface.co/h-uns>

Image: `docker pull webis/valueeval24-arthur-schopenhauer-ensemble:1.0.0`

**Team Hierocles of Alexandria [51].**<sup>10</sup> The team used both the multi-lingual dataset and English translations and incorporated sentence sequence information. More specifically, their approach predicts values for a sentence from an input text that consists of the previous two sentences concatenated with the target sentence. The two preceding sentences contained special tokens to represent any values assigned to them. During training and validation the true labels were employed, but during testing the predicted labels of the previous sentences were leveraged. The team fine-tuned different RoBERTa [4] and DeBERTa [3] models for English and XLM-RoBERTa [5] models for the multi-lingual dataset, with the best performing one being XLM-RoBERTa-xl [52]. Moreover, they developed a custom model architecture for multi-label text classification consisting of multiple classification heads. Each classification head focused on a different language for the multi-lingual dataset. The custom model architecture was adapted and employed for the English-translated dataset as well. After preliminary experiments concerning loss functions, class weights and various thresholds, they used the binary cross-entropy loss with logits as their loss function and selected an optimal classification threshold for each value. The approach is trained to tackle both subtasks 1 and 2.

**Team Philo of Alexandria [53].**<sup>11</sup> The team used the English translations of the dataset and analyzed the sentences independently. They approached subtask 1 as a multi-label problem and fine-tuned DeBERTa (deberta-base [3]) after initial experiments with several models. They employ the same base model for subtask 2 and fine-tune it to classify each text pair of sentence and human value name into either attaining or constraining.

**Team SCaLAR NITK (code name: Peter Abelard) [54].** The team used the English translations of the dataset and analyzed the sentences independently. They experimented with SVMs, KNNs, decision trees, hierarchical classification, transformer models and large language models. Based on preliminary experiments, they fine-tuned a RoBERTa [4] model for both subtasks (multi-label and binary classification, respectively).

#### 4.4. Task Evaluation

Following ValueEval'23 [2], submissions are evaluated using standard macro  $F_1$ -score over all values. The same metric is used for the new subtask 2. The submission format allowed participants to submit only one run file for both subtasks (same format as the `labels.tsv`), but the scores for the subtasks are calculated independently of each other from the same file as follows. Each submission includes for each sentence and value a confidence score (between 0 and 1) for both attained and constrained polarity. If the sum of the two numbers is above 0.5, the submission is evaluated as having predicted that the sentence refers to that value (subtask 1). For subtask 2, only the sentence-value pairs are considered for which the sentence refers to the value according to the ground-truth. For these pairs, the submission is evaluated as having predicted the attainment polarity for which it produced the larger confidence score.

Table 3 shows the results for the best-performing approaches per team for both subtasks. The best-performing approach for subtask 1 is the one of team Hierocles of Alexandria that uses XLM-RoBERTa-xl, the previous sentences, and is trained specifically for subtask 1. Overall, multilingual models performed best, with also the second-in-place employing such a model. Rarer values are overall detected worse, with the exception of the zero-shot approach by team Eric Fromm (especially Humility), indicating insufficient training data. Several teams achieved top scores for subtask 2. Overall, this binary classification task is, as once can expect, much easier than subtask 1. However, most teams clearly focused their efforts on subtask 1, so there is likely more room for improvement.

---

<sup>10</sup>Code: <https://github.com/SotirisLegkas/Touche-ValueEval24-Hierocles-of-Alexandria>

Image: `docker pull webis/valueeval24-hierocles-of-alexandria:1.0.0`

<sup>11</sup>Code: <https://github.com/VictorMYeste/touche-human-value-detection>

Models: <https://huggingface.co/VictorYeste/deberta-based-human-value-detection>

<https://huggingface.co/VictorYeste/deberta-based-human-value-stance-detection>

Image: `docker pull victoriyeste/valueeval24-philo-of-alexandria-deberta-cascading`

**Table 3**

F<sub>1</sub>-score of the best submission per team (measured by overall F<sub>1</sub>-score) on the test dataset for subtasks 1 and 2, and whether the submission used the original multilingual dataset or the automatic translation to English (EN). Baseline submissions (“Aristotle”) and the best-performing submission of ValueEval’23 (“Adam Smith,” without re-training) are shown in gray. The appendix contains tables with all submissions on page 23 and 24.

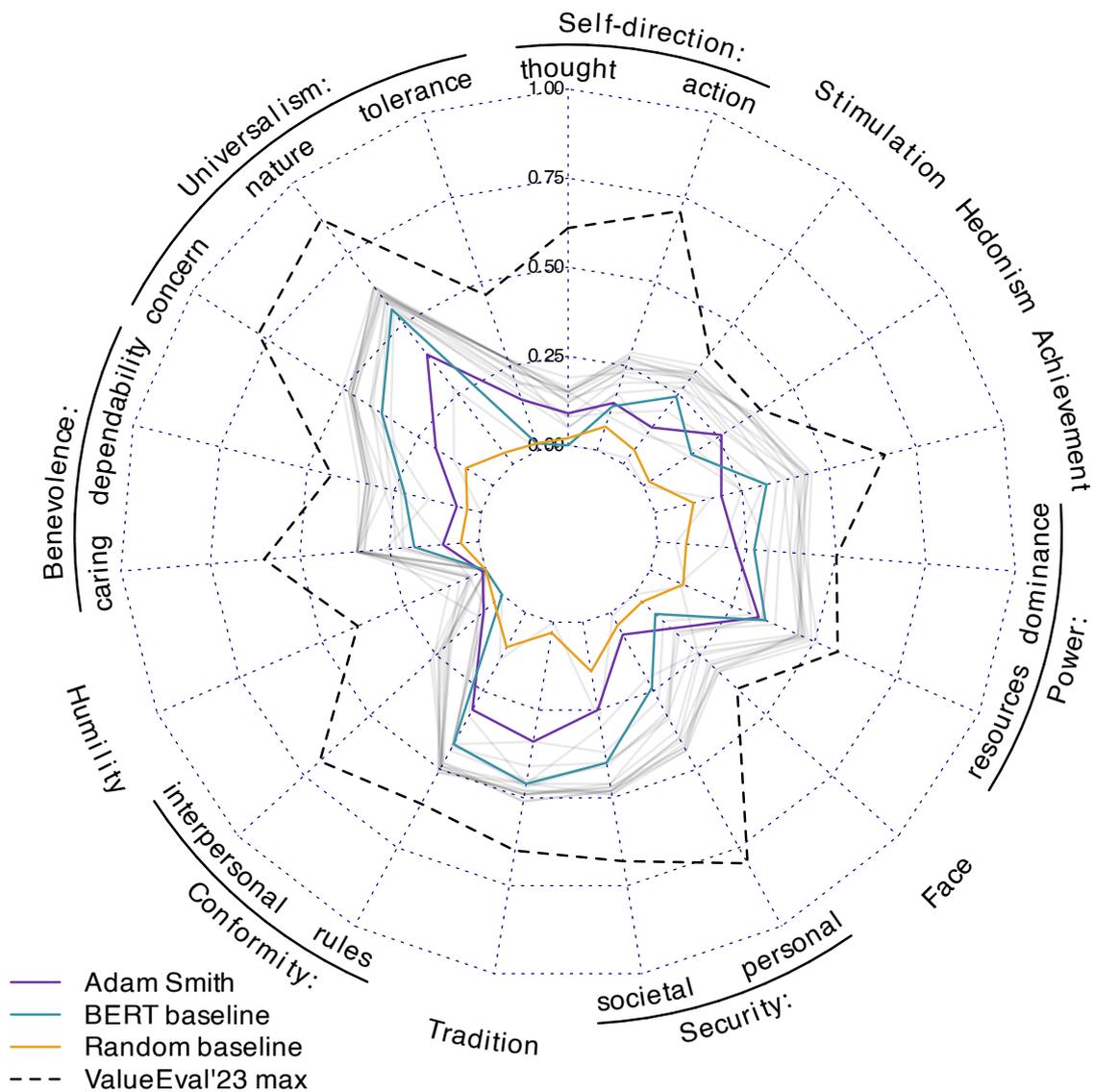
Subtask 1

Team	Lang.	F <sub>1</sub> -score																			
		Overall	Self-direction: thought	Self-direction: action	Stimulation	Hedonism	Achievement	Power: dominance	Power: resources	Face	Security: personal	Security: societal	Tradition	Conformity: rules	Conformity: interpersonal	Humility	Benevolence: caring	Benevolence: dependability	Universalism: concern	Universalism: nature	Universalism: tolerance
Hierocles of Alexandria [51]	multil.	39	15	27	30	37	45	42	49	31	42	49	46	51	24	00	34	33	47	63	27
Arthur Schopenhauer [48]	multil.	35	12	24	33	35	40	37	47	24	38	46	49	50	19	00	32	31	46	60	27
Philo of Alexandria [53]	EN	28	08	22	27	31	35	31	34	17	33	40	47	42	09	00	21	28	40	57	21
SCaLAR NITK [54]	EN	28	05	17	27	27	38	34	38	15	34	40	41	43	07	00	23	26	37	56	16
Edward Said [49]	EN	28	05	17	11	15	25	31	34	16	32	41	45	44	06	05	10	23	41	57	27
Erich Fromm [50]	EN	25	15	10	10	18	25	18	09	24	21	30	46	33	09	15	26	15	41	55	20
Lawrence Kohlberg	EN	25	08	11	19	23	31	22	31	11	28	37	34	42	09	00	21	23	34	54	18
Aristotle (BERT)	EN	24	00	13	24	16	32	27	35	08	24	40	46	42	00	00	18	22	37	55	02
Adam Smith	EN	20	09	14	13	26	19	22	33	14	07	25	34	31	07	01	10	07	19	39	15
John Shelby Spong	EN	07	00	00	02	00	16	05	11	00	01	28	00	15	00	00	00	00	13	27	00
Alain Badiou	EN	07	00	00	02	00	16	05	11	00	01	28	00	15	00	00	00	00	13	27	00
Aristotle (random)	EN	06	02	07	05	02	11	08	10	03	04	14	03	11	03	00	05	04	09	04	02

Subtask 2

Arthur Schopenhauer [48]	multil.	83	77	83	85	88	87	73	84	80	82	84	78	80	79	74	91	89	86	85	81
Edward Said [49]	EN	83	77	82	85	88	88	79	80	77	84	84	85	80	80	76	90	86	85	85	78
Philo of Alexandria [53]	EN	82	85	80	85	91	86	79	80	78	85	80	82	77	78	77	93	89	84	83	79
Aristotle (BERT)	EN	81	83	79	86	88	84	77	80	74	84	81	78	78	79	87	89	86	85	81	78
John Shelby Spong	EN	81	81	77	83	88	88	77	79	76	83	82	85	76	81	84	90	85	81	81	79
Alain Badiou	EN	81	81	77	83	88	88	77	79	76	83	82	85	76	81	84	90	85	81	81	79
Hierocles of Alexandria [51]	multil.	77	73	73	77	75	78	77	79	71	78	79	77	78	74	25	74	77	78	84	71
SCaLAR NITK [54]	EN	77	69	72	78	73	79	77	79	71	78	81	79	77	70	70	77	76	79	80	71
Erich Fromm [50]	EN	70	71	69	73	70	72	74	73	67	60	66	76	70	68	73	75	71	70	73	67
Lawrence Kohlberg	EN	66	81	77	83	80	70	76	63	56	33	45	85	63	46	84	90	79	69	70	60
Aristotle (random)	EN	52	51	47	54	52	53	55	53	52	52	50	54	53	49	45	53	56	52	49	56

To get a visual impression of the performance of the submissions, the radar plot in Figure 2 shows the F<sub>1</sub>-score of each submission for each value as lines. As the plot shows, almost all submissions have improved for all values compared to the random baseline (orange). However, all lines lie within the black dashed boundary of the maximum F<sub>1</sub>-scores achieved by last year’s submissions on last year’s dataset [55, 2]. This shows that the difficulty of the prediction task has increased compared to last year,



**Figure 2:** F<sub>1</sub>-score for each value of each submission for subtask 1, with lines corresponding to one submission each. Baselines and the best-performing submission of ValueEval’23 (“Adam Smith,” without re-training) are colored. For a comparison across years, the black dashed line shows the maximum F<sub>1</sub>-score achieved by any submission of ValueEval’23 on the ValueEval’23 main test set. The farther a line is from the center, the better the prediction for the respective value.

mainly due to the much rarer values. The difference between the datasets of the two years can also be seen in the line of the Adam Smith classifier (purple) in comparison to that of the BERT baseline (teal): Adam Smith performs worse than the specially trained BERT baseline (also: overall F<sub>1</sub>-score 0.20 vs. 0.24) since it was not retrained for the new dataset, even though in ValueEval’23 it significantly improved over the BERT baseline (0.56 vs. 0.42).

If one compares last year’s hull of submission lines (black dashed line, “ValueEval’23 max”) with this year’s equivalent hull, one sees that some values in this year’s data set are particularly difficult to predict. The visual spread between these hull lines is particularly large for the values Self-direction: thought, Self-direction: action, Security: personal, Conformity: interpersonal and Humility. A likely explanation for this is that these values are expressed very differently in the underlying source data. We therefore conclude that value detectors are not yet robust across all text genres and that further data sets in different genres are needed to achieve this goal.

## 5. Task 2: Multilingual Ideology and Power Identification in Parliamentary Debates

The study of parliamentary debates is crucial to understand the decision processes in the parliaments and their societal impacts. The goal of this task is to automatically identify two important aspects of parliamentary debates: the political orientation of the party of the speaker, and the role of the party of the speaker in the governance of the country or the region. Identifying these underlying aspects of parliamentary debates enables automated comprehension of these discussions, the decisions that these discussions lead to, and their consequences.

### 5.1. Task Definition

Both subtasks were defined as binary classification tasks: Given a parliamentary speech, (1) predict the political orientation of the party of the speaker on the *left–right* spectrum, and (2) predict whether the speaker belongs to one of the governing parties or the opposition. The first task is relatively well studied, and there have been some recent shared tasks on identifying political orientation [32, 33]. Unlike the earlier tasks, our data set includes multiple parliaments and languages, and is based on parliamentary debates. To the best of our knowledge, automatic identification of governing role–power–has not been studied earlier.

### 5.2. Data Description

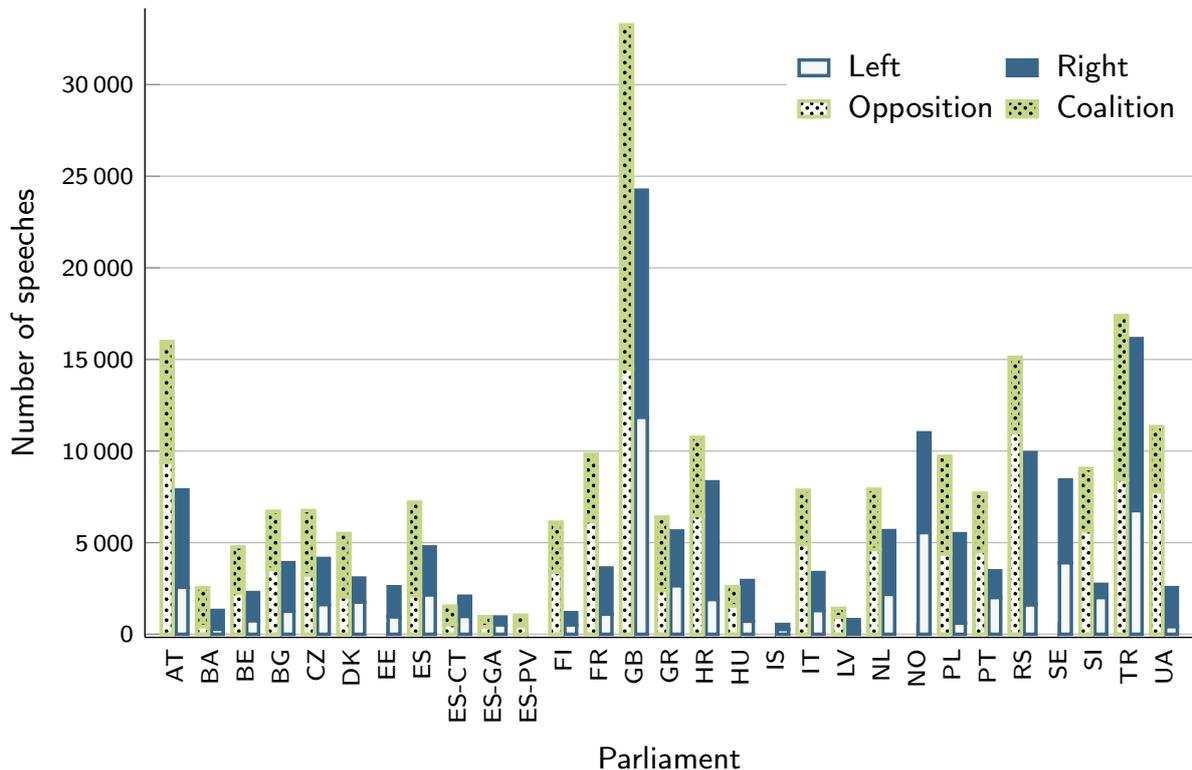
The source of the data for this task is the ParlaMint [56], a uniformly encoded and annotated corpus of transcripts of parliamentary speeches from multiple national and regional parliaments.<sup>12</sup> The transcripts are The ParlaMint version 4.0 used for the task includes data from the following national and regional parliaments: Austria (AT), Bosnia and Herzegovina (BA), Belgium (BE), Bulgaria (BG), Czechia (CZ), Denmark (DK), Estonia (EE), Spain (ES), Catalonia (ES-CT), Galicia (ES-GA), Basque Country (ES-PV), Finland (FI), France (FR), Great Britain (GB), Greece (GR), Croatia (HR), Hungary (HU), Iceland (IS), Italy (IT), Latvia (LV), The Netherlands (NL), Norway (NO), Poland (PL), Portugal (PT), Serbia (RS), Sweden (SE), Slovenia (SI), Turkey (TR) and Ukraine (UA). The labels for both subtasks are also coded in the ParlaMint corpora. For the sake of simplicity, we formulate both tasks as binary classification tasks. For both tasks, the main challenge in the creation of a dataset is to minimize the effects of covariates. Even though the instances to classify are speeches, the annotations are based on the party membership of the speaker. As a result, underlying variables like party membership, or speaker identity perfectly covary with ideology and power in most cases.

As a trade-off between data size, and for reducing the effect of covariates, we opt for a speaker-based sampling. First, to discourage, to some extent, the classifiers from relying on author identification, we sample at most 20 speeches of a single speaker. This is also important for introducing variation into the dataset, as the number of speeches from each speaker follows a power-law distribution: While a small number of speakers tend to deliver most of the speeches, e.g., party or party group leaders, most speakers have relatively few speeches. The distribution of speeches or speakers to include in training and test sets is also important for proper evaluation. For the ideology task, the set of speakers in the training and test sets are disjoint. The ideal dataset split for the power identification task requires a different constraint: training and test sets should include speeches from the same speaker with different power roles. To come as close as possible to this ideal split, we opt for a best-effort training–test split. When possible, we make sure that the speakers in the test set are also available in the training set with the opposite power role. Otherwise, we randomly sample more speakers to obtain the test set.

For evaluation, we set the test set size to 2 000 instances for both subtasks (100 to 200 speakers depending on the individual corpus and the task). Despite multiple speeches from each speaker, due to

---

<sup>12</sup>Although all transcripts are obtained through the data published by the respective parliaments, the method for obtaining the transcripts vary, such as scraping the web site of the parliament, extracting from published PDF files, and obtaining through an API provided by the parliament. For details, we refer to [56].



**Figure 3:** Overview of the Touché24 ideology and power identification dataset. The bars show the training set for both subtasks for each parliament. Test set sizes are approximately 2000 speeches for all parliaments.

missing annotations and the lack of diversity of orientation in some parliaments, the disjoint speakers constraint mentioned above results in a small number of instances in the training set for some of the parliaments. Not all parliamentary data provides both labels. Some countries do not have the opposition–governing party distinction, and for the Galician parliament, the number and distribution of orientation labels did not result in a test set that was large enough. Figure 3 shows the training set sizes for each parliament. The test set size for all parliaments is approximately 2000 speeches. We do not provide a validation set. We provide further details on the data set and the sampling procedure in a separate publication [57].<sup>13</sup>

In addition to the original speech transcripts and labels, we also provide automatic English translations, an anonymized speaker ID and the speaker’s sex in the data for both tasks. Except the speaker ID, which is not in the test sets.

Both data sets exhibit a mild class and text length imbalance between parliaments. The data set’s size was a technical challenge for some participants. The average text length is approximately 600 space-separated tokens, which is larger than the maximum accepted by many of the pretrained language models. Moreover, the data set is also large overall (more than 3GB uncompressed).

### 5.3. Participant Approaches

In 2024, 9 teams participated in this task and submitted 52 runs. We added a baseline for comparison. Unlike the ValueEval task, where pretrained language models were the dominant classifiers, for this task many participants preferred traditional, ‘computationally light’ approaches. A possible reason may be the large text size which is more costly to process with larger systems. Most teams, even the teams that used language models with large context sizes, truncated the texts to alleviate computational requirements. Some of the interesting improvements include ensemble of classifiers, data augmentation

<sup>13</sup>Training and test data are available at <https://zenodo.org/doi/10.5281/zenodo.10450640>, and <https://zenodo.org/doi/10.5281/zenodo.11061649> respectively.

through back-translation and synonym replacement, multi-task learning, additional features, such as sentiment scores, and the use of domain-specific models.

**Baselines.** We provided only a single logistic regression baseline with tf-idf weighted character n-grams. The baseline is intentionally kept simple to encourage participation by early researchers, and reduce the computation requirements.

**Team Policy Parsing Panthers [58].** The team did a set of experiments with original transcripts and their English translations, using various deep pretrained models, including BERT [47], mBERT [47], RoBERTa [4], XLM-RoBERTa [5], DeBERTa-v3 [3] Gemma [59] and ensembles of these models. This team presents an extensive set of approaches, and their analyses. A few interesting approaches worth mentioning in this short summary includes (1) Data augmentation and balancing through back-translation, (2) experiments with additional metadata, (3) multi-task learning, (4) the use of automatically obtained polarity labels, and increasing the number of instances in the training set of the orientation subtask by using the matching speaker IDs in the power dataset. This team participated in both subtasks for all parliaments.

**Team Trojan Horses [60].** The team experimented with improving the logistic regression baseline, as well as fine-tuning BERT. They used the English translations and participated in both subtasks for the majority of the parliaments.

**Team Pixel Phantoms [61].** The team experimented with some of the traditional classifiers (SVMs, logistic regression and decision trees) using the English translations provided. As well as tf-idf weighted features, they also extracted text embeddings from DistilBERT [62], through Sentence BERT [63]. They participated in both subtasks for the majority of the parliaments.

**Team Ssnites [64].** The team fine-tuned BERT for the majority of parliaments and both subtasks. They relied on the English translations provided, and participated in both subtasks for the majority of the parliaments.

**Team Hale Lab [65].** After some initial experiments with BERT, the team used a variety of classification methods including simple feed-forward networks, and LSTMs. The features for the models were either bag-of-words features weighted with tf-idf, or the multilingual LASER [66] embeddings. They used the original (untranslated) data, using various libraries for tokenization and preprocessing, and participated in both subtasks for the majority of the parliaments.

**Team Vayam Solve Kurmaha [67].** This team also experimented with multiple traditional classification methods (SVM, kNN, random forests) and their ensembles, using the English translations. The team also used data augmentation through synonym replacement. They participated in both subtasks for the majority of the parliaments.

**Team Gerber [68].** The team used a convolutional neural network (CNN) for the task without any pretrained embeddings. They used the original transcripts only, and participated in both subtasks for the majority of the parliaments.

**Team JU\_NLP\_DID [69].** The team used SVM classifiers with tf-idf features, participating in both subtasks for the majority of the parliaments. They also make use of automatic sentiment labels as an additional feature.

**Table 4**

F<sub>1</sub>-scores of all submissions on power identification task. Baseline scores are shown in gray.

Team	F <sub>1</sub> -score																												
	Overall	AT	BA	BE	BG	CZ	DK	EE	ES	ES-CT	ES-GA	FI	FR	GB	GR	HR	HU	IS	IT	LV	NL	NO	PL	PT	RS	SE	SI	TR	UA
Policy Parsing Panthers	79	77	51	71	77	63	84	64	94	80	98	77	75	92	89	65	87	71	77	67	71	82	88	95	79	95	78	93	83
gerber	63	60	45	54	62	52	56	00	77	66	76	54	58	76	72	51	69	00	60	49	59	00	72	69	64	00	58	84	73
HALE Lab	61	56	44	59	60	52	56	52	76	69	84	52	48	74	71	43	67	57	60	49	53	61	62	67	55	77	49	83	60
Pixel Phantoms	59	58	49	56	56	47	56	54	72	64	75	59	58	72	71	55	68	57	57	54	60	54	59	54	51	61	47	78	56
Ssnites	59	50	53	55	53	50	61	52	61	58	64	55	56	64	59	53	60	58	53	51	56	66	71	64	64	75	58	79	53
Trojan Horses	59	61	25	57	61	51	60	57	72	67	00	33	60	73	74	53	71	55	66	00	60	61	68	63	00	74	00	80	68
INSA Passau	59	60	53	54	61	47	57	53	63	61	66	34	58	69	59	56	66	56	56	54	56	58	69	55	61	66	51	80	62
JU_NLP_DID	57	53	42	42	55	51	60	57	69	57	70	00	50	71	63	43	60	55	61	47	56	59	51	67	48	73	46	77	57
Baseline	56	52	42	45	53	52	56	47	72	65	67	54	43	74	74	43	57	39	56	45	51	62	46	63	53	75	39	84	58

**Team INSA Passau [70].** The team also experimented with multiple approaches, where some of their submissions were focused on orientation identification and a smaller number of parliaments. The methods used included training SVMs, fine-tuning BERT-based models (pre)trained on legal documents [71, 72] and finetuning and zero- and few-shot prompting the Llama [73] version 3 models with varying sizes (which were released during while the shared task was running).

#### 5.4. Task Evaluation

We use macro-averaged F<sub>1</sub>-score as the main evaluation metric for both subtasks. Similar to the ValueEval task, the participants were encouraged to submit confidence scores, where a score over 0.5 is interpreted as class 1 and otherwise 0.

Table 4 and Table 5 present the overall best-performing approaches per team for the ideology and power subtasks respectively. The best scores for both tasks are from the team Policy Parsing Panthers. The team used an ensemble of multiple models, with multiple improvements including data augmentation and multitask learning. Results on the tables do not include approaches that were focused on only one or a small number of parliaments. A noteworthy focused submission for only GB and ideology subtask by the team INSA Passau based on fine-tuning the most recent Llama 3 model achieved the second-best result for this parliament. Although the results on both tasks are higher than the baseline we provided, the variation in the scores indicate that there is quite some room for improvement for each of the approaches.

As the results show, as formulated in this task, identifying orientation is slightly more difficult than identifying power. The overall success of the systems on a particular parliament depends on, among others, size and class distribution of the training data, and composition of the parliament. For example, there is a general trend (with some exceptions) that for parliaments with few or no government and opposition role changes in the data (e.g., HU, PL, and TR) the roles are easier to predict than for parliaments with more varied composition and more role changes( e.g., AT, BA, and UA).

## 6. Task 3: Image Retrieval/Generation for Arguments (joint task with ImageCLEF)

Images provide powerful visual communication, are usually perceived before text is read, and can appeal directly to our emotions. The goal of this task is to find images that convey premises. The proper use

**Table 5**

F<sub>1</sub>-scores of the best submissions per team (as measured by overall F<sub>1</sub>-score) on power identification task for each parliament. Baseline scores are shown in gray.

Team	F <sub>1</sub> -score																											
	Overall	AT	BA	BE	BG	CZ	DK	ES	ES-CT	ES-GA	ES-PV	FI	FR	GB	GR	HR	HU	IT	LV	NL	PL	PT	RS	SI	TR	UA		
Policy Parsing Panthers	83	88	56	74	81	78	87	88	91	98	90	80	82	83	95	75	97	78	75	74	90	85	84	81	94	65		
HALE Lab	70	69	46	61	68	69	70	65	85	88	78	65	67	75	82	68	88	69	62	64	78	65	69	61	84	49		
Trojan Horses	69	72	57	63	67	63	68	69	82	85	74	39	66	72	83	67	86	72	64	64	74	65	75	62	83	56		
gerber	68	68	51	60	66	64	63	72	80	86	74	60	71	72	68	63	87	52	63	64	77	66	73	58	84	48		
Vayam Solve Kurmaha	68	48	48	65	69	68	69	72	83	87	76	35	66	47	85	67	88	72	62	68	75	67	75	63	85	48		
Pixel Phantoms	66	70	50	59	63	65	69	65	64	77	69	61	64	73	72	57	80	69	58	62	70	66	69	60	80	52		
Baseline	64	66	45	61	68	64	56	65	78	83	71	56	66	71	63	60	86	43	51	62	76	62	65	53	83	46		
JU_NLP_DID	63	68	47	55	58	57	67	60	78	55	72	00	59	00	77	65	83	71	47	63	70	63	54	56	78	43		
INSA Passau	62	67	45	60	66	65	54	65	00	00	00	56	66	72	56	61	85	45	52	64	77	62	63	54	84	47		
Ssnites	60	66	45	58	60	61	61	62	58	62	60	60	65	60	69	65	79	62	54	57	62	58	60	57	61	46		

of an image can increase the persuasiveness of an argument. In this regard, images can increase the pathos [74], which is the effect an argument has on its audience.

## 6.1. Task Definition

This observation leads to our task, in which participants are asked to find images based on an argument that help to convey the premise of the argument. In this context, “convey” is meant in broad terms; it can represent what is described in the argument, but it can also show a generalization (e.g., a symbolic image that illustrates a related abstract concept) or a specialization (e.g., a concrete example). There is a difference between verbal language and images. Verbal language provides clear but limited information, while images provide more information than written words, but are not as precise [75]. Therefore, images alone can be ambiguous and difficult to understand without context, e.g. when they refer to symbolism. For this reason, we offer the option of submitting a rationale together with the image. The rationale is an explanatory statement that assists in understanding the picture. For example, it can be a caption or contextual information about the image. The image and the rationale are evaluated together to see how this combination conveys the premise. Participants can choose to use a retrieval approach, where they submit images from a provided dataset, or a generation-based approach, where suitable images can be generated using a model of their choice. In each submission, a participant can submit up to 10 images in a ranking order for an argument.

## 6.2. Data Description

For the task we prepared a dataset<sup>14</sup> containing 136 arguments and over 9000 images. The arguments were generated with GPT-4 [76] and correspond to 24 topics. The topics were taken from various IBM datasets<sup>15</sup> and previous Touché Shared Tasks<sup>16</sup>. Each generated argument consists of a premise and a claim, and can take a pro or con stance on the topic. An example of an argument can be seen in Fig. 4. Each of the images in the dataset is tagged with additional information, such as the URL and content of the corresponding website. In addition, we have provided an analysis of each image using the Google

<sup>14</sup><https://zenodo.org/records/11045831>

<sup>15</sup>[https://research.ibm.com/haifa/dept/vst/debating\\_data.shtml](https://research.ibm.com/haifa/dept/vst/debating_data.shtml)

<sup>16</sup><https://touche.webis.de/shared-tasks.html>

```

<argument>
  <id>36062-a-3</id>
  <topic>Should boxing be banned?</topic>
  <premise>
    The idea of winning through intentional infliction of pain and harm
    to another person can nurture a violent and destructive mentality.
  </premise>
  <claim>
    Boxing poses both physical and psychological threats to
    participants, hence it should be banned.
  </claim>
  <stance>pro</stance>
  <type>ANECDOTAL</type>
</argument>

```

**Figure 4:** Example argument from the data set. The argument consists of an id, a premise and a claim. The data also indicates the topic of the argument, as well as the argument’s stance on the topic. The type element indicates that the arguments relies on anecdotal evidence. Only arguments of this type are used in our dataset.

<p><b>Topic: Should boxing be banned?</b></p> <hr/> <p><b>Premise:</b> The idea of winning through intentional infliction of pain and harm to another person can nurture a violent and destructive mentality.</p> <hr/> <p><b>Claim:</b> Boxing poses both physical and psychological threats to participants, hence it should be banned.</p> <hr/> <p><b>Rationale:</b> The image captures a boxing match in progress, with two men standing in the ring. One of them is wearing a red glove and appears to be getting hit by his opponent’s punch. The other man is also wearing a red glove, likely as part of his attire for the match. The boxers are focused on their performance, with one of them holding his mouth open while taking a blow from his opponent. The scene showcases the intensity and determination of these athletes during the competition.</p>	
---	---

**Figure 5:** Example of a submission on the topic that boxing should be banned. The image conveys the premise, as it clearly shows the threats of boxing for the participants. In this example, the participating team has not opted for its own rationale. In such cases, the automatically generated image caption with LLaVA is used as the default rationale. Image taken from <https://qph.cf2.quoracdn.net/main-qimg-d962ebc7a954e8b1a5ec8bae6bde6662-lq>

Cloud Vision API, as well as an automatically generated caption using LLaVA [77]. An example of a submission can be seen in Figure 5.

### 6.3. Participant Approaches

In 2024, 2 teams participated in this task and submitted 8 runs. All teams chose the retrieval-approach. Moreover, we added 2 baseline runs for comparison.

**Baselines** The first baseline is BM25, where the corresponding documents are the image captions from the data set and the query is the premise of the argument. In the second baseline, keywords are first extracted from the image captions. Then embeddings for the premise of an argument and the keywords are generated with SBERT [63]. A corresponding relevance score is calculated based on the cosine similarity between the embeddings and averaging them. The most relevant images are selected for submission.

**DS@GT [78].** The team uses CLIP [6] to embed each argument and each image in a common embedding space. The first approach ranks images by cosine similarity of the embeddings. The second approach compares for each argument the 40 highest ranked images to images that are generated to support or attack the argument. The most similar images are submitted. For image generation, Stable Diffusion v2-1 [79] was used.

**HTW-DIL [80].** The team has chosen an approach inspired by DPR [7]. It applies a fine-tuned multimodal Moondream model based on the Phi 1.5 LLM [81] and uses SigLIP [82] for its vision capabilities. To generate synthetic training data, the team uses GPT-4 to generate arguments from the available image/web page data. Combinations of positive and negative argument-image pairs are used for training. The results are obtained by maximising the cosine similarity for argument and image embeddings. To enable comparability, the team also adopted a standard approach of embedding the corresponding website content of the images and each argument using OpenAI’s Ada model<sup>17</sup> and selecting the most similar pairs.

## 6.4. Task Evaluation

For each argument and each submission, the best 5 images together with the rationales are evaluated by a human expert. This expert knows neither the rank of the image nor the team that submitted it. To facilitate the annotation, we prepared a narrative for each argument that describes what a conveying image should generally show. Therefore, each combination of image, argument and rationale is rated on a three-point Likert scale from 0 to 2, where 0 means that the image does not convey the premise at all, 1 stands for partial conveyance and 2 means that the image conveys the premise completely. A total of 5,061 image, argument and rationale triples were annotated. For seven topics, only very few relevant images could be submitted by the participating teams, so we removed these topics, resulting in a total number of 104 arguments for the final evaluation. For each submission, we first calculated the NDCG score for each argument. For the required IDCG, we have considered all submitted image, argument and rationale triples submitted for the corresponding argument. The final score of a submission is the average of all NDCG scores for all arguments.

Table 6 shows the results for both teams and baselines. For all three NDCG measures, team HTW-DIL achieved the highest scores with the submission that ranks the results using OpenAI’s Ada embeddings of the website content and the argument—thus not using the image at all. Other submissions were similar to the top-performing submissions from previous years. As such an approach was not successful in earlier years, likely this year’s updated task description, which provides complete arguments instead of mere topics, enabled the top performance of this approach. The performance of combined approaches is yet to be tested. And as the achieved scores below 0.5 show, the identification of images that convey a specific argument is still a very challenging task.

## 7. Conclusion

The fifth edition of the Touché lab on argumentation systems featured three tasks: (1) Human Value Detection, (2) Ideology and Power Identification in Parliamentary Debates, and (3) Image Retrieval/-Generation for Arguments. In contrast to previous years, the focus this year was more on classification than retrieval tasks. Furthermore, two of the three tasks were multilingual, although automatic English transcriptions were provided to facilitate participation. We expanded the scope of Touché with the new tasks on human values and political power and orientation. In addition, we methodically extended the retrieval task by allowing participants to generate images instead of retrieving them. Unfortunately, no team submitted generated images in the end.

Of the 68 registered teams, 20 participated in the tasks and submitted a total of 81 runs. Participants mainly used classification architectures, with BERT and variants still very dominant, although more clas-

---

<sup>17</sup><https://platform.openai.com/docs/models/embeddings>

**Table 6**

NDCG values for the participating teams and their approaches for the top 5, top 3 and most relevant image(s). The approaches are sorted according to the NDCG@5 value. The winning approach Ada-Summary from HTW-Dil refers to OpenAI Ada embeddings. The various Moondream models use Moondream embeddings for image or web content text or for both together. The suffix EP indicates the number of training epochs. Base-CLIP uses CLIP embeddings and Generated-Image-CLIP uses CLIP in combination with images generated by Stable Diffusion.

Rank	Team	Approach	NDCG@5	NDCG@3	NDCG@1
1	HTW-DIL	Ada-Summary	0.428	0.409	0.404
2	HTW-DIL	Moondream-Text	0.363	0.355	0.356
3	HTW-DIL	Moondream-Default-Image-Text	0.293	0.302	0.317
4	Baseline	BM25	0.284	0.273	0.293
5	Baseline	SBERT	0.232	0.225	0.221
6	DS@GT	Generated-Image-CLIP	0.180	0.178	0.197
7	HTW-DIL	Moondream-Image-Text-EP3	0.150	0.163	0.183
8	HTW-DIL	Moondream-Image	0.146	0.155	0.178
9	DS@GT	Base-CLIP	0.123	0.111	0.106
10	HTW-DIL	Moondream-Image-Text-EP2	0.120	0.140	0.178

sical machine learning models were also used in the Ideology and Power Identification in Parliamentary Debates task. Generative models, on the other hand, were rarely used. The Image Retrieval/Generation for Arguments task changed to seeking images for a specific argument rather than a topic, and the best-performing submission used an approach that was not successful for the previous task definitions: it ranked images by the embedding similarity between the argument and the web page that contains the image—and thus ignored the actual image content.

We plan to continue Touché as a collaborative platform for researchers in argumentation systems. All Touché resources are freely available, including topics, manual relevance, argument quality, and stance judgments, and submitted runs from participating teams. These resources and other events such as workshops will help to further foster the community working on argumentation systems.

## Acknowledgments

This work was partially supported by the European Commission under grant agreement GA 101070014 (<https://openwebsearch.eu>) and the German Research Foundation under project 455911521 (LARGA) as part of the SPP 1999 (RATIO). The ideology and power identification shared task has been supported by CLARIN ERIC, under the ParlaMint project (<https://www.clarin.eu/parlamint>).

## References

- [1] J. Kiesel, Ç. Çöltekin, M. Heinrich, M. Fröbe, M. Alshomary, B. D. Longueville, T. Erjavec, N. Handke, M. Kopp, N. Ljubešić, K. Meden, N. Mirzakhmedova, V. Morkevičius, T. Reitis-Münstermann, M. Scharfbillig, N. Stefanovitch, H. Wachsmuth, M. Potthast, B. Stein, Overview of Touché 2024: Argumentation Systems, in: L. Goeuriot, P. Mulhem, G. Quénot, D. Schwab, L. Soulier, G. M. D. Nunzio, P. Galuščáková, A. G. S. de Herrera, G. Faggioli, N. Ferro (Eds.), *Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Fifteenth International Conference of the CLEF Association (CLEF 2024)*, Lecture Notes in Computer Science, Springer, Berlin Heidelberg New York, 2024.
- [2] J. Kiesel, M. Alshomary, N. Mirzakhmedova, M. Heinrich, N. Handke, H. Wachsmuth, B. Stein, SemEval-2023 Task 4: ValueEval: Identification of Human Values behind Arguments, in: R. Kumar,

- A. K. Ojha, A. S. Dođruöz, G. D. S. Martino, H. T. Madabushi (Eds.), Proc. of SemEval, ACL, 2023, pp. 2287–2303. doi:10.18653/v1/2023.semeval-1.313.
- [3] P. He, X. Liu, J. Gao, W. Chen, DeBERTa: decoding-enhanced BERT with disentangled attention, in: Proc. of ICLR, 2021. URL: <https://openreview.net/forum?id=XPZlaotutsD>.
- [4] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, RoBERTa: A Robustly Optimized BERT Pretraining Approach, CoRR (2019). URL: <http://arxiv.org/abs/1907.11692>.
- [5] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, V. Stoyanov, Unsupervised Cross-lingual Representation Learning at Scale, in: D. Jurafsky, J. Chai, N. Schluter, J. R. Tetreault (Eds.), Proc. of ACL, ACL, 2020, pp. 8440–8451. doi:10.18653/v1/2020.acl-main.747.
- [6] A. Radford, et al., Learning Transferable Visual Models From Natural Language Supervision, in: M. Meila, T. Zhang (Eds.), Proc. of ICML, volume 139, PMLR, 2021, pp. 8748–8763. URL: <https://proceedings.mlr.press/v139/radford21a.html>.
- [7] V. Karpukhin, et al., Dense Passage Retrieval for Open-Domain Question Answering, in: Proc. of EMNLP, ACL, 2020, pp. 6769–6781. doi:10.18653/v1/2020.emnlp-main.550.
- [8] S. H. Schwartz, Are There Universal Aspects in the Structure and Contents of Human Values?, Journal of Social Issues (1994) 19–45. doi:10.1111/j.1540-4560.1994.tb01196.x.
- [9] T. Bench-Capon, Persuasion in Practical Argument Using Value-based Argumentation Frameworks, Journal of Logic and Computation 13 (2003) 429–448. doi:10.1093/logcom/13.3.429.
- [10] S. H. Schwartz, J. Cieciuch, M. Vecchione, E. Davidov, R. Fischer, C. Beierlein, A. Ramos, M. Verkasalo, J.-E. Lönnqvist, K. Demirutku, et al., Refining the Theory of Basic Individual Values, Journal of personality and social psychology (2012). doi:10.1037/a0029393.
- [11] M. Scharfbillig, V. Ponizovskiy, Z. Pasztor, J. Keimer, G. Tirone, Monitoring Social Values in Online Media Articles on Child Vaccinations, Technical Report, European Commission’s Joint Research Centre, Luxembourg, 2022. doi:10.2760/86884.
- [12] M. Scharfbillig, L. Smillie, D. Mair, M. Sienkiewicz, J. Keimer, R. Pinho Dos Santos, H. Vinagreiro Alves, E. Vecchione, L. Scheunemann, Values and Identities - a Policymaker’s Guide, Technical Report, European Commission’s Joint Research Centre, Luxembourg, 2021. doi:10.2760/349527.
- [13] P. Ammanabrolu, L. Jiang, M. Sap, H. Hajishirzi, Y. Choi, Aligning to Social Norms and Values in Interactive Narratives, in: M. Carpuat, M. de Marneffe, I. V. M. Ruíz (Eds.), Proc. of NAACL-HLT 2022, ACL, 2022, pp. 5994–6017. doi:10.18653/v1/2022.naacl-main.439.
- [14] R. Liu, C. Jia, G. Zhang, Z. Zhuang, T. X. Liu, S. Vosoughi, Second Thoughts are Best: Learning to Re-Align With Human Values from Text Edits, Advances in Neural Information Processing Systems 35 (2022) 181–196.
- [15] L. Qiu, Y. Zhao, J. Li, P. Lu, B. Peng, J. Gao, S. Zhu, ValueNet: A New Dataset for Human Value Driven Dialogue System, in: Proc. of AAAI, AAAI Press, 2022, pp. 11183–11191. doi:10.1609/aaai.v36i10.21368.
- [16] M. Forbes, J. D. Hwang, V. Shwartz, M. Sap, Y. Choi, Social Chemistry 101: Learning to Reason about Social and Moral Norms, in: B. Webber, T. Cohn, Y. He, Y. Liu (Eds.), Proc. of EMNLP, ACL, 2020, pp. 653–670. doi:10.18653/v1/2020.emnlp-main.48.
- [17] J. Kiesel, M. Alshomary, N. Handke, X. Cai, H. Wachsmuth, B. Stein, Identifying the Human Values behind Arguments, in: S. Muresan, P. Nakov, A. Villavicencio (Eds.), Proc. of ACL, ACL, 2022, pp. 4459–4471. doi:10.18653/v1/2022.acl-long.306.
- [18] D. Fišer, J. Lenardič, CLARIN resources for parliamentary discourse research, in: D. Fišer, M. Eskevich, F. de Jong (Eds.), Proc. of LREC, ELRA, 2018.
- [19] J. Lenardič, D. Fišer, CLARIN Resource Families: Parliamentary Corpora, 2023. <https://www.clarin.eu/resource-families/parliamentary-corpora>, accessed on 2024-07-09.
- [20] G. Glavaš, F. Nanni, S. P. Ponzetto, Computational Analysis of Political Texts: Bridging Research Efforts Across Communities, in: 57th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts, ACL, 2019, pp. 18–23. doi:10.18653/v1/P19-4004.
- [21] G. Abercrombie, R. Batista-Navarro, Sentiment and position-taking analysis of parliamentary

- debates: a systematic literature review, *Journal of Computational Social Science* 3 (2020) 245–270.
- [22] A. Arian, M. Shamir, The primarily political functions of the left-right continuum, *Comparative politics* 15 (1983) 139–158.
- [23] F. Vegetti, D. Širinić, Left–right categorization and perceptions of party ideologies, *Political Behavior* 41 (2019) 257–280.
- [24] T. van Dijk, *Discourse and Power*, Bloomsbury Publishing, 2008.
- [25] N. Fairclough, *Critical Discourse Analysis: The Critical Study of Language*, Longman applied linguistics, Taylor & Francis, 2013. doi:10.4324/9781315834368.
- [26] N. Fairclough, *Language and Power, Language In Social Life*, Taylor & Francis, 2013. doi:10.4324/9781315838250.
- [27] M. D. Conover, B. Gonçalves, J. Ratkiewicz, A. Flammini, F. Menczer, Predicting the political alignment of Twitter users, in: *Proc. of PASSAT and SocialCom, IEEE*, 2011, pp. 192–199. doi:10.1109/PASSAT/SocialCom.2011.34.
- [28] S. Gerrish, D. M. Blei, Predicting Legislative Roll Calls from Text, in: L. Getoor, T. Scheffer (Eds.), *Proc. of ICML*, Omnipress, 2011, pp. 489–496.
- [29] D. Preoțiuc-Pietro, Y. Liu, D. Hopkins, L. Ungar, Beyond Binary Labels: Political Ideology Prediction of Twitter Users, in: R. Barzilay, M.-Y. Kan (Eds.), *Proc. of ACL, ACL*, 2017, pp. 729–740. doi:10.18653/v1/P17-1068.
- [30] F. Pla, L.-F. Hurtado, Political Tendency Identification in Twitter using Sentiment Analysis Techniques, in: J. Tsujii, J. Hajic (Eds.), *Proc. of Coling, Dublin City University and ACL*, 2014, pp. 183–192. URL: <https://aclanthology.org/C14-1019>.
- [31] C. Chen, D. Walker, V. Saligrama, Ideology Prediction from Scarce and Biased Supervision: Learn to Disregard the “What” and Focus on the “How”!, in: A. Rogers, J. Boyd-Graber, N. Okazaki (Eds.), *Proc. of ACL (Volume 1: Long Papers), ACL, Toronto, Canada*, 2023, pp. 9529–9549. doi:10.18653/v1/2023.acl-long.530.
- [32] J. A. García-Díaz, et al., Overview of PoliticES 2022: Spanish Author Profiling for Political Ideology, *Procesamiento del Lenguaje Natural* 69 (2022) 265–272. doi:10.26342/2022-69-23.
- [33] D. Russo, et al., PoliticIT at EVALITA 2023: Overview of the political ideology detection in Italian texts task, in: *Proc. of EVALITA, volume 3473 of CEUR Workshop Proceedings, CEUR-WS.org*, 2023. URL: <https://ceur-ws.org/Vol-3473/paper7.pdf>.
- [34] G. M. Kurtoğlu Eskişar, Ç. Çöltekin, Emotions Running High? A Synopsis of the state of Turkish Politics through the ParlaMint Corpus, in: D. Fišer, M. Eskevich, J. Lenardič, F. de Jong (Eds.), *Proc. of ParlaCLARIN, ELRA*, 2022, pp. 61–70. URL: <https://aclanthology.org/2022.parlaclarin-1.10>.
- [35] M. Mochtak, P. Rupnik, N. Ljubešić, The ParlaSent Multilingual Training Dataset for Sentiment Identification in Parliamentary Proceedings, in: N. Calzolari, M.-Y. Kan, V. Hoste, A. Lenci, S. Sakti, N. Xue (Eds.), *Proc. of LREC, ELRA and ICCL*, 2024, pp. 16024–16036. URL: <https://aclanthology.org/2024.lrec-main.1393>.
- [36] O. Tarkka, J. Koljonen, M. Korhonen, J. Laine, K. Martiskainen, K. Elo, V. Laippala, Automated Emotion Annotation of Finnish Parliamentary Speeches Using GPT-4, in: D. Fiser, M. Eskevich, D. Bordon (Eds.), *Proc. of ParlaCLARIN, ELRA and ICCL*, 2024, pp. 70–76. URL: <https://aclanthology.org/2024.parlaclarin-1.11>.
- [37] C. Navarretta, D. Haltrup Hansen, Government and opposition in Danish parliamentary debates, in: D. Fiser, M. Eskevich, D. Bordon (Eds.), *Proc. of ParlaCLARIN, ELRA and ICCL*, 2024, pp. 154–162. URL: <https://aclanthology.org/2024.parlaclarin-1.23>.
- [38] I. J. Dove, On images as evidence and arguments, in: F. H. van Eemeren, B. Garssen (Eds.), *Topical Themes in Argumentation Theory: Twenty Exploratory Studies*, Argumentation Library, Springer Netherlands, Dordrecht, 2012, pp. 223–238. doi:10.1007/978-94-007-4041-9\_15.
- [39] F. Dunaway, Images, emotions, politics, *Modern American History* 1 (2018) 369–376. doi:10.1017/mah.2018.17.
- [40] G. Roque, Visual argumentation: A further reappraisal, in: F. H. van Eemeren, B. Garssen (Eds.), *Topical Themes in Argumentation Theory, volume 22*, Springer Netherlands, 2012, pp. 273–288. doi:10.1007/978-94-007-4041-9\_18.

- [41] I. Grancea, Types of visual arguments, *Argumentum. Journal of the Seminar of Discursive Logic, Argumentation Theory and Rhetoric* 15 (2017) 16–34.
- [42] D. Dimitrov, B. Bin Ali, S. Shaar, F. Alam, F. Silvestri, H. Firooz, P. Nakov, G. Da San Martino, SemEval-2021 Task 6: Detection of Persuasion Techniques in Texts and Images, in: *Proc. of SemEval, ACL*, 2021, pp. 70–98. URL: <https://aclanthology.org/2021.semeval-1.7>. doi:10.18653/v1/2021.semeval-1.7.
- [43] H. Wachsmuth, N. Naderi, Y. Hou, Y. Bilu, V. Prabhakaran, T. A. Thijm, G. Hirst, B. Stein, Computational argumentation quality assessment in natural language, in: *Proc. of EACL*, 2017, pp. 176–187. URL: <https://aclanthology.org/E17-1017>.
- [44] M. Fröbe, M. Wiegmann, N. Kolyada, B. Grahm, T. Elstner, F. Loebe, M. Hagen, B. Stein, M. Potthast, Continuous Integration for Reproducible Shared Tasks with TIRA.io, in: J. Kamps, L. Goeuriot, F. Crestani, M. Maistro, H. Joho, B. Davis, C. Gurrin, U. Kruschwitz, A. Caputo (Eds.), *Proc. of ECIR, Lecture Notes in Computer Science*, Springer, 2023, pp. 236–241. doi:10.1007/978-3-031-28241-6\_20.
- [45] N. Stefanovitch, J. Piskorski, Holistic Inter-Annotator Agreement and Corpus Coherence Estimation in a Large-scale Multilingual Annotation Campaign, in: H. Bouamor, J. Pino, K. Bali (Eds.), *Proc. of EMNLP, ACL*, 2023, pp. 71–86. doi:10.18653/v1/2023.emnlp-main.6.
- [46] M. V. Nguyen, V. D. Lai, A. P. B. Veyseh, T. H. Nguyen, Trankit: A Light-Weight Transformer-based Toolkit for Multilingual Natural Language Processing, in: D. Gkatzia, D. Seddah (Eds.), *Proc. of EACL, ACL*, 2021, pp. 80–90. doi:10.18653/v1/2021.eacl-demos.10.
- [47] J. Devlin, M. Chang, K. Lee, K. Toutanova, BERT: pre-training of deep bidirectional transformers for language understanding, in: J. Burstein, C. Doran, T. Solorio (Eds.), *Proc. of NAACL-HLT, ACL*, 2019, pp. 4171–4186. doi:10.18653/v1/N19-1423.
- [48] H. Yunis, Arthur Schopenhauer at Touché 2024: Multi-Lingual Text Classification Using Ensembles of Large Language Models, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), *Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org*, 2024.
- [49] A. Aydin, S. Shaar, C. Cardie, Edward said at touché: Human values classification, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), *Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org*, 2024.
- [50] M. Morren, R. Mishra, Eric fromm at touché: Prompts vs finetuning, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), *Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org*, 2024.
- [51] S. Legkas, C. Christodoulou, M. Zidianakis, D. Koutrintzes, G. Petasis, M. Dagioglou, Hierocles of alexandria at touché: Multi-task & multi-head custom architecture with transformer-based models for human value detection, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), *Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org*, 2024.
- [52] N. Goyal, J. Du, M. Ott, G. Anantharaman, A. Conneau, Larger-Scale Transformers for Multilingual Masked Language Modeling, in: A. Rogers, I. Calixto, I. Vulic, N. Saphra, N. Kassner, O. Camburu, T. Bansal, V. Shwartz (Eds.), *Proc. of RepL4NLP@ACL-IJCNLP, ACL*, 2021, pp. 29–33. doi:10.18653/v1/2021.REPL4NLP-1.4.
- [53] V. Yeste, M. C. Ardanuy, P. Rosso, Philo of Alexandria at Touché: A Cascade Model Approach to Human Value Detection, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), *Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org*, 2024.
- [54] P. K. D. K. C. Reddy, A. M, ScaLAR NITK at Touché: Human Value Detection, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), *Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org*, 2024.
- [55] N. Mirzakhmedova, J. Kiesel, M. Alshomary, M. Heinrich, N. Handke, X. Cai, V. Barriere, D. Dastgheib, O. Ghahroodi, M. Sadraei, E. Asgari, L. Kawaletz, H. Wachsmuth, B. Stein, The Touché23-ValueEval Dataset for Identifying Human Values behind Arguments, in: N. Calzolari, M.-Y. Kan,

- V. Hoste, A. Lenci, S. Sakti, N. Xue (Eds.), Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), International Committee on Computational Linguistics, 2024.
- [56] T. Erjavec, M. Ogrodniczuk, et al., The ParlaMint corpora of parliamentary proceedings, LREC 57 (2022) 415–448. doi:10.1007/s10579-021-09574-0.
- [57] Ç. Çöltekin, M. Kopp, M. Katja, V. Morkevicius, N. Ljubešić, T. Erjavec, Multilingual Power and Ideology identification in the Parliament: a reference dataset and simple baselines, in: D. Fiser, M. Eskevich, D. Bordon (Eds.), 4th Workshop on Creating, Analysing, and Increasing Accessibility of Parliamentary Corpora, ELRA and ICCL, 2024, pp. 94–100. URL: <https://aclanthology.org/2024.parlaclarin-1.14>.
- [58] O. Palmqvist, J. Jiremalm, P. Picazo-Sanchez, Policy Parsing Panthers at Touché: Ideology and Power Identification in Parliamentary Debates, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [59] T. Mesnard, et al., Gemma: Open Models Based on Gemini Research and Technology, 2024. doi:10.48550/arXiv.2403.08295. arXiv:2403.08295.
- [60] P. Mirunalini, A. Koushik, D. S. D. Seshan, Trojan Horses at Touché: Logistic regression for classification of political debates, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [61] J. Hariharakrishnan, J. S. P. Mirunalini, Pixel Phantoms at Touché: Ideology and power identification in parliamentary debates using linear SVC, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [62] V. Sanh, L. Debut, J. Chaumond, T. Wolf, DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter, 2020. arXiv:1910.01108.
- [63] N. Reimers, I. Gurevych, Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks, in: Proc. of EMNLP, ACL, 2019, pp. 3982–3992. doi:10.18653/v1/D19-1410.
- [64] K. V. K. S. K. A. M. P. S. N. Ssnites at Touché: Ideology and power identification in parliamentary debates using BERT model, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [65] S. Sevitha, M. Patel, S. Shevgoor, Team Hale Lab at Touché 2024: Ideology and power identification in parliamentary debates, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [66] M. Artetxe, H. Schwenk, Massively Multilingual Sentence Embeddings for Zero-Shot Cross-Lingual Transfer and Beyond, Transactions of the Association for Computational Linguistics 7 (2019) 597–610. doi:10.1162/tacl\_a\_00288.
- [67] S. Shwetha, S. Kamath, S. Balaji, S. N. S. R. S. Narayanan, Vayam Solve Kurmaha at Touché: Power Identification in Parliamentary Speeches Using TFIDF Vectorizer and SVM Classifier, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [68] C. Gerber, gerber at Touché: Ideology and power identification in parliamentary debates 2024, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [69] A. Khurshid, D. Das, R. Khaskel, S. Datta, JU\_NLP\_DID at Touché, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [70] M. Andruszak, A. Alhamzeh, E. Egyed-Zsigmond, A. Carlsson, J. Leydet, Y. Otiefy, Team INSA Passau at Touché: Multi-lingual parliamentary speech classification, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the

- Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
- [71] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, N. Aletras, I. Androutsopoulos, LEGAL-BERT: The muppets straight out of law school, in: T. Cohn, Y. He, Y. Liu (Eds.), Findings of ACL: EMNLP 2020, ACL, 2020, pp. 2898–2904. doi:10.18653/v1/2020.findings-emnlp.261.
  - [72] L. Zheng, N. Guha, B. R. Anderson, P. Henderson, D. E. Ho, When does pretraining help?: assessing self-supervised learning for law and the CaseHOLD dataset of 53,000+ legal holdings, in: Proc. of ICAIL, ACM, 2021, pp. 159–168. doi:10.1145/3462757.3466088.
  - [73] H. Touvron, et al., LLaMA: Open and Efficient Foundation Language Models, 2023. doi:10.48550/arxiv.2302.13971.
  - [74] C. Rapp, Aristotle’s Rhetoric, in: E. N. Zalta, U. Nodelman (Eds.), The Stanford Encyclopedia of Philosophy, Metaphysics Research Lab, Stanford University, 2023.
  - [75] J. E. Kjeldsen, Virtues of visual argumentation: How pictures make the importance and strength of an argument salient, 2013.
  - [76] J. Achiam, et al., GPT-4 Technical Report, 2024. arXiv:2303.08774.
  - [77] H. Liu, C. Li, Q. Wu, Y. J. Lee, Visual instruction tuning, 2023. arXiv:2304.08485.
  - [78] B. Ostrower, P. Aphiwetsa, Ds@gt at touché: Image search and ranking via clip and image generation, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
  - [79] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, in: Proc. of CVPR, IEEE, 2022, pp. 10674–10685. doi:10.1109/CVPR52688.2022.01042.
  - [80] T. Janusko, A. Kämpf, D. Keiling, J. Knick, D. S. M. Thiele, Htw-dil at touché: Multimodal dense information retrieval for arguments, in: G. Faggioli, N. Ferro, P. Galuščáková, A. G. S. de Herrera (Eds.), Working Notes of the Conference and Labs of the Evaluation Forum (CLEF 2024), CEUR Workshop Proceedings, CEUR-WS.org, 2024.
  - [81] Y. Li, S. Bubeck, R. Eldan, A. D. Giorno, S. Gunasekar, Y. T. Lee, Textbooks Are All You Need II: phi-1.5 technical report, 2023. URL: <https://arxiv.org/abs/2309.05463>.
  - [82] X. Zhai, B. Mustafa, A. Kolesnikov, L. Beyler, Sigmoid loss for language image pre-training, in: Proc. of ICCV, IEEE Computer Society, 2023, pp. 11941–11952. doi:10.1109/iccv51070.2023.01100.

## A. Extended Results

**Table 7**

Achieved  $F_1$ -score of all submissions on the test dataset for subtasks 1, and whether the submission used the original multilingual dataset or the automatic translation to English (EN). Baseline submissions (“Aristotle”) and the winning submission of ValueEval’23 (“Adam Smith,” without re-training) are shown in gray.

Team	Approach	Lang.	F <sub>1</sub> -score																			
			Overall	Self-direction: thought	Self-direction: action	Stimulation	Hedonism	Achievement	Power: dominance	Power: resources	Face	Security: personal	Security: societal	Tradition	Conformity: rules	Conformity: interpersonal	Humility	Benevolence: caring	Benevolence: dependability	Universalism: concern	Universalism: nature	Universalism: tolerance
Hierocles of Alexandria	XLM-RoBERTa-xl tkns 38 train+val	multil.	39	15	27	30	37	45	42	49	31	42	49	46	51	24	00	34	33	47	63	27
Hierocles of Alexandria	XLM-RoBERTa-xl tkns 19 train+val	multil.	38	15	27	31	36	43	41	51	32	44	49	48	51	23	00	34	35	50	63	24
Hierocles of Alexandria	RoBERTa-l tkns weighted-19 train	EN	37	19	23	31	32	40	41	45	31	43	48	51	48	26	11	34	33	48	60	27
Hierocles of Alexandria	RoBERTa-l tkns 19 train+val	EN	37	16	28	33	35	43	38	48	28	44	48	51	49	27	05	34	27	48	61	27
Hierocles of Alexandria	DeBERTa-v2-xxl tkns 19 train	EN	37	15	26	32	32	44	40	45	32	41	47	49	50	24	05	34	33	48	62	27
Hierocles of Alexandria	RoBERTa-l tkns 38 train+val	EN	37	12	24	32	36	42	39	46	28	43	47	49	49	22	00	34	32	47	61	27
Hierocles of Alexandria	XLM-RoBERTa-l 19 train+val	multil.	36	15	28	35	35	44	39	47	28	40	48	49	50	20	08	33	32	47	60	24
Hierocles of Alexandria	XLM-RoBERTa-l tkns 19 train+val	EN	35	14	25	30	28	41	40	46	25	40	48	48	48	20	05	34	30	46	59	25
Arthur Schopenhauer		multil.	35	12	24	33	35	40	37	47	24	38	46	49	50	19	00	32	31	46	60	27
Hierocles of Alexandria	XLM-RoBERTa-l tkns weighted-19 train	multil.	34	13	20	28	28	37	37	45	22	33	46	46	49	21	04	32	32	47	63	21
Philo of Alexandria		EN	28	08	22	27	31	35	31	34	17	33	40	47	42	09	00	21	28	40	57	21
SCaLAR NITK		EN	28	05	17	27	27	38	34	38	15	34	40	41	43	07	00	23	26	37	56	16
Edward Said		EN	28	05	17	11	15	25	31	34	16	32	41	45	44	06	05	10	23	41	57	27
Erich Fromm		EN	25	15	10	10	18	25	18	09	24	21	30	46	33	09	15	26	15	41	55	20
Lawrence Kohlberg		EN	25	08	11	19	23	31	22	31	11	28	37	34	42	09	00	21	23	34	54	18
Aristotle	BERT baseline	EN	24	00	13	24	16	32	27	35	08	24	40	46	42	00	00	18	22	37	55	02
Adam Smith	ValueEval’23	EN	20	09	14	13	26	19	22	33	14	07	25	34	31	07	01	10	07	19	39	15
John Shelby Spong		EN	07	00	00	02	00	16	05	11	00	01	28	00	15	00	00	00	00	13	27	00
Alain Badiou	1	EN	07	00	00	02	00	16	05	11	00	01	28	00	15	00	00	00	00	13	27	00
Alain Badiou	2	EN	07	00	00	02	00	16	05	11	00	01	28	00	15	00	00	00	00	13	27	00
Aristotle	Random baseline	EN	06	02	07	05	02	11	08	10	03	04	14	03	11	03	00	05	04	09	04	02

**Table 8**

Achieved  $F_1$ -score of all submissions on the test dataset for subtasks 2, and whether the submission used the original multilingual dataset or the automatic translation to English (EN). Baseline submissions (“Aristotle”) are shown in gray.

Team	Approach	Lang.	F <sub>1</sub> -score																			
			Overall	Self-direction: thought	Self-direction: action	Stimulation	Hedonism	Achievement	Power: dominance	Power: resources	Face	Security: personal	Security: societal	Tradition	Conformity: rules	Conformity: interpersonal	Humility	Benevolence: caring	Benevolence: dependability	Universalism: concern	Universalism: nature	Universalism: tolerance
Arthur Schopenhauer		multil.	83	77	83	85	88	87	73	84	80	82	84	78	80	79	74	91	89	86	85	81
Edward Said		EN	83	77	82	85	88	88	79	80	77	84	84	85	80	80	76	90	86	85	85	78
Philo of Alexandria		EN	82	85	80	85	91	86	79	80	78	85	80	82	77	78	77	93	89	84	83	79
Aristotle	BERT baseline	EN	81	83	79	86	88	84	77	80	74	84	81	78	78	79	87	89	86	85	81	78
John Shelby Spong		EN	81	81	77	83	88	88	77	79	76	83	82	85	76	81	84	90	85	81	81	79
Alain Badiou	1	EN	81	81	77	83	88	88	77	79	76	83	82	85	76	81	84	90	85	81	81	79
Alain Badiou	2	EN	81	81	77	83	88	88	77	79	76	83	82	85	76	81	84	90	85	81	81	79
Hierocles of Alexandria	XLM-RoBERTa-xl tkns 38 train+val	multil.	77	73	73	77	75	78	77	79	71	78	79	77	78	74	25	74	77	78	84	71
Hierocles of Alexandria	RoBERTa-l tkns 38 train+val	EN	77	72	72	78	74	78	78	78	73	78	78	78	77	73	22	78	77	78	82	74
SCaLAR NITK		EN	77	69	72	78	73	79	77	79	71	78	81	79	77	70	70	77	76	79	80	71
Erich Fromm		EN	70	71	69	73	70	72	74	73	67	60	66	76	70	68	73	75	71	70	73	67
Hierocles of Alexandria	RoBERTa-l tkns weighted-19 train	EN	69	73	71	76	71	71	74	68	65	55	59	77	69	64	70	77	74	70	72	66
Hierocles of Alexandria	XLM RoBERTa-l tkns weighted-19 train	multil.	69	71	70	74	72	71	74	68	66	61	57	76	69	65	70	76	75	70	73	65
Hierocles of Alexandria	RoBERTa-l tkns 19 train+val	EN	68	73	71	77	72	71	74	67	63	55	59	78	69	62	70	77	72	68	72	65
Hierocles of Alexandria	XLM-RoBERTa-l tkns 19 train+val	EN	68	72	72	77	74	72	74	67	64	56	55	77	68	64	70	77	76	69	73	63
Hierocles of Alexandria	XLM-RoBERTa-l 19 train+val	multil.	68	72	72	77	73	72	74	67	64	57	57	77	68	63	70	75	73	69	72	64
Hierocles of Alexandria	XLM-RoBERTa-xl tkns 19 train+val	multil.	68	71	71	76	72	72	74	66	64	56	56	79	68	62	70	75	74	69	72	66
Hierocles of Alexandria	DeBERTa-v2-xxl tkns 19 train	EN	68	71	70	77	74	72	76	66	63	57	57	77	69	61	70	78	75	68	73	65
Lawrence Kohlberg		EN	66	81	77	83	80	70	76	63	56	33	45	85	63	46	84	90	79	69	70	60
Aristotle	Random baseline	EN	52	51	47	54	52	53	55	53	52	50	54	53	49	45	53	56	52	49	56	56

**Table 9**

F<sub>1</sub>-scores of all submissions on ideology identification task. Baseline scores are shown in gray.

		F <sub>1</sub> -score																												
Team	Approach	Overall	AT	BA	BE	BG	CZ	DK	EE	ES	ES-CT	ES-GA	FI	FR	GB	GR	HR	HU	IS	IT	LV	NL	NO	PL	PT	RS	SE	SI	TR	UA
Policy Parsing Panthers	Regular model predictions	79	77	52	70	75	59	82	64	91	79	98	76	75	93	87	66	87	71	76	66	71	82	87	92	77	95	76	93	82
Policy Parsing Panthers	Dataset vulnerabilities	79	77	51	71	77	63	84	64	94	80	98	77	75	92	89	65	87	71	77	67	71	82	88	95	79	95	78	93	83
INSA Passau	l370b0sl	79													79															
Policy Parsing Panthers	immature-havarti	78	78	53	69	74	60	80	65	92	78	98	76	74	92	88	67	86	70	74	68	70	81	88	91	75	94	76	92	83
Policy Parsing Panthers	plain-bugle	78	78	53	69	74	60	80	65	92	78	98	76	74	92	88	67	86	70	74	68	70	81	88	91	75	94	76	92	83
INSA Passau	l370b2sl2	78													78															
INSA Passau	l370b2sl1	77													77															
Policy Parsing Panthers	wicker-fowl	76	74	51	68	73	56	79	62	88	73	96	70	72	91	82	65	83	67	71	67	67	79	86	89	74	90	73	91	79
INSA Passau	l370b0sl_v2	76													76															
Policy Parsing Panthers	dense-loop	73	72	53	66	72	57	66	62	85	74	95	72	69	89	79	60	81	67	65	64	64	71	83	81	74	87	73	90	82
INSA Passau	l38b0sl	73													73															
INSA Passau	llama_ft	68													68															
gerber	constant-feta	63	60	45	54	62	52	56		77	66	76	54	58	76	72	51	69		60	49	59		72	69	64		58	84	73
HALE Lab	universal-triangle	61	56	44	59	60	52	56	52	76	69	84	52	48	74	71	43	67	57	60	49	53	61	62	67	55	77	49	83	60
INSA Passau	l38b2sl1	61													61															
INSA Passau	bert_all_lang	59	60	53	54	61	47	57	53	63	61	66	34	58	69	59	56	66	56	56	54	56	58	69	55	61	66	51	80	62
Pixel Phantoms	balanced-photon	59	58	49	56	56	47	56	54	72	64	75	59	58	72	71	55	68	57	57	54	60	54	59	54	51	61	47	78	56
Pixel Phantoms	collinear-cuisine	59	57	49	56	56	48	56	54	72	65	75	59	58	72	72	54	69	57	57	55	59	53	59	54	51	60	47	78	56
Pixel Phantoms	run1	59	57	49	56	56	48	56	54	72	65	75	59	58	72	72	54	69	57	57	55	59	53	59	54	51	60	47	78	56
Pixel Phantoms	run2	59	58	49	56	56	47	56	54	72	64	75	59	58	72	71	55	68	57	57	54	60	54	59	54	51	61	47	78	56
Ssnites	foggy-destination	59	50	53	55	53	50	61	52	61	58	64	55	56	64	59	53	60	58	53	51	56	66	71	64	64	75	58	79	53
Trojan Horses	convoluted-sonar	59	61	25	57	61	51	60	57	72	67		33	60	73	74	53	71	55	66		60	61	68	63		74		80	68
JU_NLP_DID	Ideology SVM	57	53	42	42	55	51	60	57	69	57	70		50	71	63	43	60	55	61	47	56	59	51	67	48	73	46	77	57
Baseline	-	56	52	42	45	53	52	56	47	72	65	67	54	43	74	74	43	57	39	56	45	51	62	46	63	53	75	39	84	58
INSA Passau	l38b2sl2	56													56															
INSA Passau	bert_basic	55	53	43	43	58	49	53	51	59	53	57	29	46	70	51	51	57	53	51	45	55	57	66	58	64	68	50	80	62
INSA Passau	svm	54	48	50	60	60	51	57	28	64	69	78	55	48	60	58	44	74	44	48	52	44	30	67	63	48	32	53	73	68
INSA Passau	logreg	54	53	42	57	54	52	55	44	72			53	46	74	74	44	55	43	51	45	53	47	45	63	51	58	41	83	60
INSA Passau	l370b_voting	52													52															
INSA Passau	zscore	51	57	47	50	39	47	38	47	52	49	59		55	67	55	52	62	48	42	48	39	44	62	59	57	60	39	55	56
Pixel Phantoms	oriented-soda	50	42	43	50	45	48	50	48	52	52	55	49	56	61	48	46	55	56	45	50	53	54	53	46	49	50	39	58	57
INSA Passau	l38b_voting	42													42															
INSA Passau	z_bert_vot	37	35	29	30	39	33	25	35	40	31	38		31	47	35	30	39	35	34	31	37	38	44	39	43	46	33	54	42
HALE Lab	crimson-highlight	37								32	34	29	35	42	36	30	42	43	38	35	45	38	37	45	31	46	34	27	35	46

