

Extracting Cultural Commonsense Knowledge at Scale

Tuan-Phong Nguyen
Max Planck Institute for Informatics
Saarbrücken, Germany
tuanphong@mpi-inf.mpg.de

Aparna Varde
Montclair State University
Montclair, New Jersey, USA
vardea@montclair.edu

Simon Razniewski
Max Planck Institute for Informatics
Saarbrücken, Germany
srazniew@mpi-inf.mpg.de

Gerhard Weikum
Max Planck Institute for Informatics
Saarbrücken, Germany
weikum@mpi-inf.mpg.de

ABSTRACT

Structured knowledge is important for many AI applications. Commonsense knowledge, which is crucial for robust human-centric AI, is covered by a small number of structured knowledge projects. However, they lack knowledge about human traits and behaviors conditioned on socio-cultural contexts, which is crucial for situative AI. This paper presents CANDLE, an end-to-end methodology for extracting high-quality cultural commonsense knowledge (CCSK) at scale. CANDLE extracts CCSK assertions from a huge web corpus and organizes them into coherent clusters, for 3 domains of subjects (geography, religion, occupation) and several cultural facets (food, drinks, clothing, traditions, rituals, behaviors). CANDLE includes judicious techniques for classification-based filtering and scoring of interestingness. Experimental evaluations show the superiority of the CANDLE CCSK collection over prior works, and an extrinsic use case demonstrates the benefits of CCSK for the GPT-3 language model. Code and data can be accessed at <https://cultural-csk.herokuapp.com/>.

1 INTRODUCTION

Motivation. Structured knowledge, often stored in knowledge graphs (KGs) [12, 39], is a key asset for many AI applications, including search, question answering, and conversational bots. KGs cover factual knowledge about notable entities such as singers, songs, cities, sports teams, etc. However, even large-scale KGs deployed in practice hardly touch on the dimension of commonsense knowledge (CSK): properties of everyday objects, behaviors of humans, and more. Some projects, such as ConceptNet [36], Atomic [32], and Ascent++ [21] have compiled large sets of CSK assertions, but are solely focused on “universal CSK”: assertions that are agreed upon by almost all people and are thus viewed as “globally true”. What is missing, though, is that CSK must often be viewed in the *context of specific social or cultural groups*: the world view of a European teenager does not necessarily agree with those of an American business person or a Far-East-Asian middle-aged factory worker.

This paper addresses this gap, by automatically compiling CSK that is conditioned on socio-cultural contexts. We refer to this as *cultural CSK* or *CCSK* for short. For example, our CCSK collection contains assertions such as:

- SUBJECT:EAST ASIA, FACET:FOOD, TOFU IS A MAJOR INGREDIENT IN MANY EAST ASIAN CUISINES, OR
- SUBJECT:FIREFIGHTER, FACET:BEHAVIOR, FIREFIGHTERS USE LADDERS TO REACH FIRES.

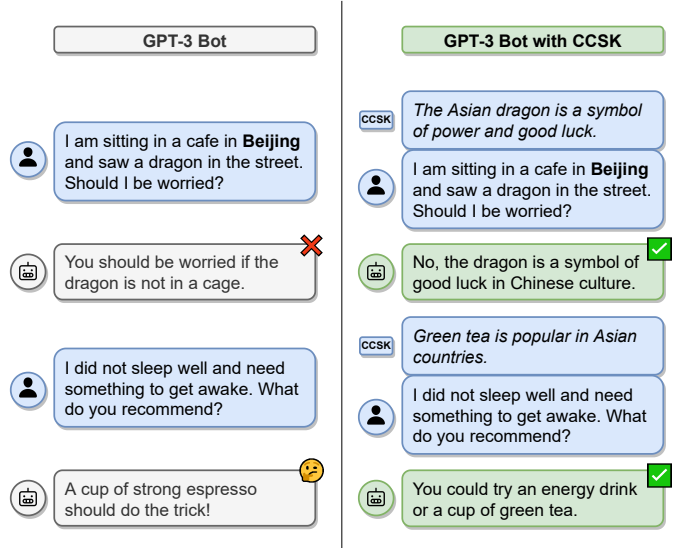


Figure 1: Human-bot conversations without and with CCSK.

The value of having a KG with this information lies in making AI applications more situative and more robust.

Consider the conversation between a human and the GPT-3 chatbot¹ shown in Fig. 1. The GPT-3-based bot, leveraging its huge language model, performs eloquently in this conversation, but completely misses the point that the user is in China, where dragons are viewed positively and espresso is difficult to get. If we prime the bot with CCSK about Far-East-Asian culture, then GPT-3 is enabled to provide culturally situative replies. If primed with CCSK about European views (not shown in Fig. 1), the bot points out that dragons are portrayed as evil monsters but do not exist in reality and recommends a strong cup of coffee.

State of the art. Mainstream KGs do not cover CCSK at all, and major CSK collections like ConceptNet contain only very few culturally contextualized assertions. To the best of our knowledge, the only prior works with data that have specifically addressed the socio-cultural dimension are the projects Quasimodo [30], StereoKG [7], and the work of Acharya et al. [1]. The latter merely contains a few hundred assertions from crowdsourcing, StereoKG uses a specialized way of automatically extracting stereotypes from QA

¹Executed at beta.openai.com/playground using the *davinci-002* model at $temp=0.7$.

forums and is still small in size, and Quasimodo covers a wide mix of general CSK and a small fraction of culturally relevant assertions. These are the three baselines to which we compare our results.

Language models (LMs) such as BERT [8] or GPT-3 [5] are another form of machine-based CSK, including CCSK, in principle. However, all LM knowledge is in latent form, captured in learned values of billions of parameters. Knowledge cannot be made explicit; we observe it only implicitly through the LM-based outputs in applications. The example of Fig. 1 demonstrates that even large LMs like GPT-3 do not perform well when socio-cultural context matters.

Approach. CCSK is expressed in text form on web pages and social media, but this is often very noisy and difficult to extract. We devised an end-to-end methodology and system, called CANDLE (Extracting Cultural Commonsense Knowledge at Scale), to automatically extract and systematically organize a large collection of CCSK assertions. For scale, we tap into the C4 web crawl [27], a huge collection of web pages. This provides an opportunity to construct a sizable CCSK collection, but also a challenge in terms of scale and noise.

The output of CANDLE is a set of 1.1M CCSK assertions, organized into 60K coherent clusters. The set is organized by 3 domains of interest – geography, religion, occupation – with a total of 386 instances, referred to as *subjects* (or cultural groups). Per subject, the assertions cover 5 *facets* of culture: food, drinks, clothing, rituals, traditions (for geography and religion) or behaviors (for occupations). In addition, we also annotate each assertion with its salient *concepts*. Examples for the computed CCSK are shown in Fig. 2.

CANDLE operates in 6 steps. First and second, we identify candidate assertions using simple techniques for *subject detection* (named entity recognition - NER, and string matching), and *generic rule-based filtering*. Third, we *classify assertions* into specific cultural facets, which is challenging because we have several combinations of cultural groups and cultural facets, making it very expensive to create specialized training data. Instead, we creatively leverage LMs pre-trained on the Natural Language Inference (NLI) task to perform zero-shot classification on our data, with judicious techniques to enhance the accuracy. Fourth we use state-of-the-art techniques for *assertion clustering*, and fifth a simple but effective method to *extract concepts* in assertions. Lastly, we combine several features to *score* the interestingness of assertions, such as frequency, specificity, distinctiveness. This way, we steer away from overly generic assertions (which LMs like GPT-3 tend to generate) and favor assertions that set their subjects apart from others.

Contributions. The main contributions of this work are:

- (1) An end-to-end methodology to extract high-quality CCSK from very large text corpora.
- (2) New techniques for judiciously classifying and filtering CCSK-relevant text snippets, and for scoring assertions by their interestingness.
- (3) A large collection of CCSK assertions for 386 subjects covering 3 domains (geography, religion, occupation) and several facets (food, drinks, clothing, traditions, rituals, behaviors).

Experimental evaluations show that the assertions in CANDLE are of significantly higher quality than those from prior works. An extrinsic use case demonstrates that our CCSK can improve

geography>country	Germany	drinks
German beer festivals in October are a celebration of beer drinking.		
geography>region	East Asia	food
Tofu is a major ingredient in many East Asian cuisines.		
geography>region	South Asia	traditions
In South Asia, henna is often used in bridal makeup or to celebrate festivals.		
occupation	lawyer	clothing
Lawyers wear suits to look professional.		
occupation	firefighter	behaviors
Firefighters run into burning buildings to save lives .		

Figure 2: Example assertions of CANDLE, with subjects (cultural groups) of cultural domains, facets and concepts.

performance of GPT-3 in question answering. Code and data can be accessed at <https://cultural-csk.herokuapp.com/>.

2 RELATED WORK

Commonsense knowledge acquisition. There is a long tradition of CSK acquisition in AI (e.g., [10, 15, 19, 28, 34]). Earlier projects, e.g., Cyc [15] and ConceptNet [19], construct commonsense knowledge graphs (CSKGs) based on large-scale human annotations. Crowdsourcing CSKG construction has been revived in the ATOMIC project [13, 32]. CSK extraction from texts has been researched in WebChild [37], TupleKB [6], Quasimodo [30], ASER [46, 47], TransOMCS [45], GenericsKB [3], and Ascent [21, 22]. Meanwhile, Ilievski et al. [14] consolidate CSK from 7 different resources into one integrated KG. Those projects, however, have their main focus on either concept-centered knowledge (e.g., ELEPHANTS HAVE TRUNKS), social interactions (e.g., X HATES Y’S GUTS, AS A RESULT, X WANTS TO YELL AT Y), or event-centered knowledge (e.g., X DRINKING COFFEE HAPPENS AFTER X POURING THE COFFEE INTO A MUG) and do not cover much cultural knowledge. Our approach also starts from texts, but focuses on cultural commonsense knowledge (CCSK), with particular challenges in knowledge representation, assertion filtering and consolidation.

Cultural commonsense knowledge. A few works have focused specifically on CCSK. An early approach by Anacleto et al. [2] gathers CSK from users from different cultures, entered via the Open Mind Common Sense portal. However, the work is limited to a few eating habits (time for meals, what do people eat in each meal?, food for party/Christmas) in 3 countries (Brazil, Mexico, USA), and without published data. Acharya et al. [1] embark on a similar manual effort towards building a cultural CSKG, limited to a few predefined predicates and answers from Amazon MTurk workers from USA and India. Shwartz [33] maps time expressions in 27 different languages to specific hours in the day, also using MTurk annotations. StereoKG [7] mines cultural stereotypes of 5 nationalities and 5 religion groups from Twitter and Reddit questions posted by their users, however, being without proper filtering, the method results in quite many noisy and inappropriate assertions. GeoMLAMA [42] defines 16 geo-diverse commonsense concepts (e.g., traffic rules, date formats, shower time) and use crowdsourcing to collect knowledge for 5 different countries in 5 corresponding languages. The dataset was used to probe multilingual pretrained language models, however,

is not shared. Moving to computer vision, Liu et al. [18] and Yin et al. [43] expand existing visual question answering datasets with images from different cultures rather than the Western world. As a result, models trained on images from the old datasets (mostly images from Western cultures) perform poorly on the newly added images. Our methodology is the first to utilize large text corpora, and it can extract CCSK in the form of natural-language sentences, for a wide range of cultural groups and facets.

Pre-trained language models and commonsense knowledge. Remarkable advances in NLP have been achieved with pre-trained language models (LMs) such as BERT [8] and GPT variants [5, 26]. LAMA [25] designs methodology and datasets to probe masked LMs in order to acquire CSK that the models implicitly store. COMET [4] is a method that finetunes autoregressive LMs on CSK triples, and it can generate possible objects for a given pair of subject-predicate. However, the quality of the generated assertions is often considerably lower than that of the training data [20]. More recently, West et al. [40] introduce a prompting technique to collect CSK by feeding GPT-3 [5] with a few human-verified CSK triples and ask it to generate new assertions. Although it was shown that the generated resource, called AutoTOMIC, is of encouraging quality, knowledge bases from LMs are inherently problematic, because there is no apparent way to trace assertions to specific sources, e.g., to understand assertion context, or to apply filters at document level.

In this work, we leverage pre-trained LMs as sub-modules in our system to help with cultural facet classification and assertion clustering. We also show that our method can produce more distinctive CCSK assertions than querying GPT-3 with prompts.

3 CCSK REPRESENTATION

Our representation of CCSK is based on the notions of *subjects* (from 3 major domains: geography, religion and occupation) and *facets*. These are the key labels for CCSK *assertions*, which are informative sentences with salient *concepts* marked up.

We assume two sets to be given:

- \mathcal{S} : A set of **subjects** (cultural groups) s_1, \dots, s_n from a cultural **domain**, e.g., based on geo-locations (United States, China, Middle East, California), religious groups (Christians, Muslims, Buddhists) or occupations (taxi driver, professor, web developer);
- \mathcal{F} : A set F_1, \dots, F_m of **facets** of culture, e.g., food, drinks, clothing, traditions, rituals, behaviors.

Note that the cultural facets need not be mutually exclusive, e.g., food assertions sometimes overlap with traditions.

Our objective is to collect a set of CCSK assertions for a given subject and facet. Existing commonsense resources store assertions in triple format (e.g., ConceptNet [36], Quasimodo [30]), semantic frames (Ascent [22]) or generic sentences (GenericsKB [3]). Although the traditional triple-based and frame-based data models are convenient for structured querying, and well suited for regular assertions like birth dates, citizenships, etc., they often falls short of capturing nuanced natural language assertions, as essential for CSK. Moreover, recent advances in pre-trained language models have made it easier to feed downstream tasks with less structured knowledge.

With CANDLE, we thus follow the approach of GenericsKB [3], and use natural-language sentences to represent assertions.

In principle, an assertion could comprise even several sentences. The longer the assertions are, however, the harder it is to discern their core. In this work, for higher precision and simplicity of computations, we only consider single sentences.

DEFINITION 1 (CULTURAL COMMONSENSE KNOWLEDGE ASSERTION). *Given a subject s and a facet F , a CCSK assertion is a triple $(s, F, sent)$ where $sent$ is a natural-language sentence about facet F of subject s .*

Since natural language often allows to express similar assertions in many different ways, and web harvesting naturally leads to discovering similar assertions multiple times, we employ clustering as an essential component in our approach.

A **cluster** (*cls*) of CCSK assertions for one subject and cultural facet contains assertions with similar meaning, and for presentation purposes, is summarized by a single summary sentence. Each cluster also comes with a score denoting its interestingness.

To further organize assertions, we also identify salient **concepts**, i.e., important terms inside assertions, that can be used for concept-centric browsing of assertion sets.

Several examples of CCSK assertions produced by CANDLE are shown in Fig. 2.

4 METHODOLOGY

We propose an end-to-end system, called CANDLE, to extract and organize CCSK assertions based on the proposed CCSK representation. Notably, our system does not require annotating new training data, but only leverages pre-trained models with judicious techniques to enhance the accuracy. The system takes in three inputs:

- an English text corpus (e.g., a large web crawl);
- a set of *subjects* (cultural groups);
- a set of *facets* of culture.

CANDLE consists of 6 modules (see Fig. 3). Throughout the system, step by step, we reduce a large input corpus (which could contain billions of documents, mostly noisy) into high-quality clusters of CCSK assertions for the given subjects and facets. Each cluster in the output is also accompanied by a representative sentence and an interestingness score. We next elaborate on each module.

4.1 Subject detection

We start the extraction by searching for sentences that contain mentions of the given subjects. These will be the candidate sentences used in the subsequent modules. To achieve high recall, we utilize generous approaches such as string matching and named entity recognition (NER), and use more advanced filtering techniques in later modules, to ensure high precision.

For the geography and religion domains, in which subjects are named entities, we use spaCy’s NER module to detect subjects. Specifically, geo-locations are detected with the GPE tag (geopolitical entities), and religions are detected with the NORP tag (nationalities or religious or political groups). For each subject, we also utilize a list of aliases for string matching, which can be the location’s alternate names (e.g., United States, the U.S., the States), or demonyms (e.g., Colombians, Chinese, New Yorker), or names for

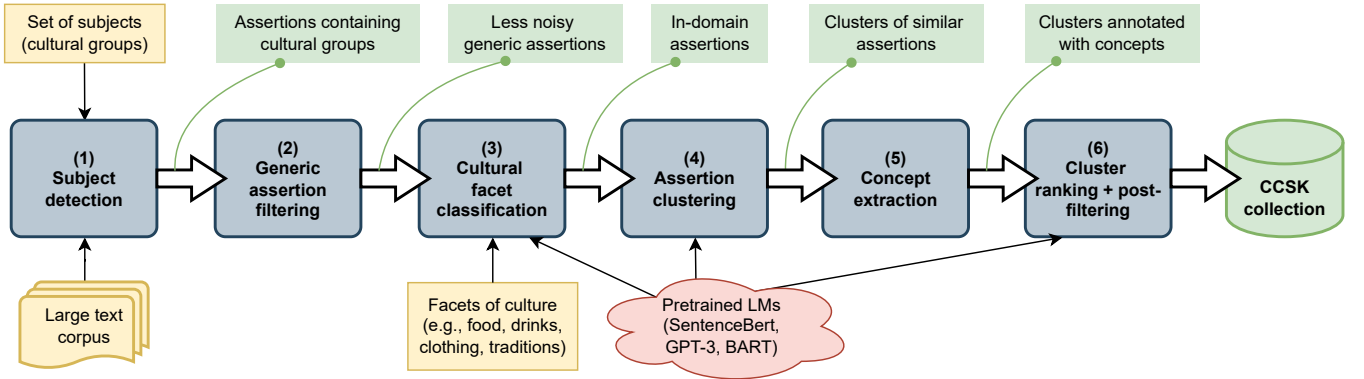


Figure 3: Architecture of CANDLE.

religious adherents (e.g., Christians, Buddhists, Muslims) - which can be detected with the NORP tag as well.

For the occupation domain, we simply use exact-phrase matching to detect candidates. Each occupation subject is enriched with its alternate names and its plural form to enhance coverage.

4.2 Generic assertion filtering

CSK aims at covering *generic assertions*, not episodic or personal experiences. For example, GERMANS LIKE THEIR CURRYWURST is a generic assertion, but I VISITED GERMANY TO EAT CURRYWURST OF THIS RESTAURANT SERVES GERMAN CURRYWURST are not.

GenericsKB [3] is arguably the most popular work on automatically identifying generic sentences in texts and it uses a set of 27 hand-crafted lexico-syntactic rules. CANDLE adopts those rules in this module. However, for each domain and facet, we adaptively drop some of the rules if they would reject valuable assertions. More details on the adaptations can be found in Appx. A.

4.3 Cultural facet classification

To organize CCSK and filter out irrelevant assertions, we classify candidate sentences into several facets of culture. Traditional methods for this classification task would require a substantial amount of annotated data to train a supervised model. The costs of data annotation are often a critical bottleneck in large-scale settings. In CANDLE, we aim to minimize the degree of human supervision by leveraging pre-trained models for zero-shot classification.

A family of pre-trained models that is suitable for our setting is *textual entailment* (a.k.a *natural language inference* - NLI): given two sentences, does one entail the other (or are they contradictory or unrelated)? Our approach to adopting such a model for cultural facet classification is inspired by the zero-shot inference method of Yin et al. [44]. Given a sentence *sent* and a facet *F*, we construct the NLI test as follows:

Input: $Premise \leftarrow sent, Hypothesis \leftarrow$ “This text is about *F*”
Output: $P[sent \in F] \leftarrow P[Premise \Rightarrow Hypothesis]$

The probability of *Premise* entailing *Hypothesis* will be taken as the probability of *sent* being labeled as *F*, denoted as $P[sent \in$

F]. For example, with sentence “German October festivals are a celebration of beer and fun”, the candidate entailments will be “This text is about drinks”, “... about food”, “... about traditions”, and so on. Multiple of these facets may yield high scores in these NLI tests.

To enhance precision, we introduce a set of *counter-labels* for topics that are completely outside the scope of CCSK, for example, politics or business. A sentence *sent* will be accepted as a good candidate for facet *F* if

$$\begin{cases} P[sent \in F] \geq \rho_+, \text{ and} \\ P[sent \in \tilde{F}] \leq \rho_- \text{ for all counter-labels } \tilde{F} \end{cases} \quad (1)$$

where ρ_+ and ρ_- are hyperparameters in the range $[0, 1]$, giving us the flexibility to tune for either precision or recall.

In our experiments, we use the BART model [16] finetuned on the MultiNLI dataset [41] for NLI tests². Our crowdsourcing evaluations show that the zero-shot classifiers with the enhanced techniques achieved high precision (see Appx. C.2).

4.4 Assertion clustering

The same assertion can be expressed in many ways in natural language. For example, FRIED RICE IS A POPULAR CHINESE DISH can also be written as FRIED RICE IS A FAMOUS DISH FROM CHINA or ONE OF THE MOST POPULAR CHINESE FOOD IS FRIED RICE. Clustering is used to group such assertions, which reduces redundancies, and allows to obtain frequency signals on assertions.

We leverage a state-of-the-art sentence embeddings method, SentenceBert [29], to compute vector representations for all assertions and use the Hierarchical Agglomerative Clustering (HAC) algorithm for clustering. Clustering is performed on assertions of each subject-facet pair.

Cluster summarization. Since each cluster can have from a few to hundreds of sentences, it is important to identify what those sentences convey, in a concise way.

One way to compute a representative assertion for a cluster is to compute the centroid of the cluster, then take its closest assertion as the representative. Yet for natural-language data, this does not work particularly well.

²Model available at <https://huggingface.co/facebook/bart-large-mnli>

In CANDLE, we therefore approach cluster summarization as a generative task, based on a state-of-the-art language model, GPT-3 [5] (see Appx. E for prompt template). Annotator-based evaluations show that GPT-generated representatives received significantly better scores than the base sentences in the clusters (see Sec. 6.1).

4.5 Concept extraction

While the cultural groups are regarded as subjects, *concepts* are akin to objects of the assertions. Identifying these concepts enables concept-focused browsing (e.g., browsing Japan assertions only about the Miso soup, etc.).

We postulate that main concepts of an assertion cluster are terms shared by many members: We extract all n -grams ($n = 1..3$) of all assertions in a cluster (excluding subjects themselves, and stop words); and retain the ones that occur in more than 60% of the assertions. If both a phrase and its sub-phrase appear, we only keep the longer phrase in the final output. Noun-phrase concepts are normalized by singularization.

4.6 Cluster ranking and post-filtering

Ranking commonsense assertions is a crucial task. Unlike encyclopedic knowledge, which is normally either true or false, precision of CSK is usually not a binary concept, as it generalizes over many groups. With CANDLE, we aim to pull out the most interesting assertions for each subject, and avoid overly generic assertions such as CHINESE FOOD IS GOOD or FIREFIGHTERS WORK HARD, which are very common in the texts.

Extracting and clustering assertions from large corpora gives us an important signal of an assertion, its *frequency*. However, ranking based on frequency alone may lead to reporting bias. As we compile a CCSK collection at large scale, it also enables us to compute the *distinctiveness* of an assertion against others in the collection. The notion of these 2 metrics can be thought of as term frequency and inverse document frequency in the established TF-IDF technique for IR document ranking [35]. Besides *frequency* and *distinctiveness*, we score the interestingness of assertion clusters based on 2 other custom metrics: *specificity* (how many objects are mentioned in the assertion?) and *domain relevance* (how relevant is the assertion to the cultural facet?).

Frequency. For each subject-facet pair, we normalize cluster sizes into the range [0, 1], using min-max normalization.

Distinctiveness. We compute the IDF of a cluster cls as follows:

$$IDF(cls) = \frac{\sum_{cls' \in CLS} size(cls')}{\sum_{cls' \in CLS} size(cls') \times \sigma(cls, cls')} \quad (2)$$

where CLS is the set of all clusters for a given facet (e.g., food) and domain (e.g., geography>country), and

$$\sigma(cls, cls') = \begin{cases} 1 & \text{if } sim(cls, cls') \geq \theta \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

Here, $sim(cls, cls')$ is the semantic similarity between the two clusters cls and cls' , and θ is a predefined threshold. In CANDLE, to reduce computation, we approximate $sim(cls, cls')$ as the similarity between their summary sentences, which can be computed as the

Table 1: Statistics of the CANDLE CCSK collection (#A: number of assertions, #C: number of clusters).

Facet	Geography		Religions		Occupations	
	#A	#C	#A	#C	#A	#C
Food	240,459	12,981	9,750	680	9,837	511
Drinks	95,394	5,923	3,079	218	3,321	227
Clothing	14,170	1,237	1,695	141	4,367	278
Rituals	116,839	8,007	74,651	3,026	22,581	1,253
Traditions	214,931	13,606	68,202	2,798	-	-
Behaviors	-	-	-	-	25,152	1,495
Other	-	-	60,483	2,292	159,239	5,461
All	681,793	41,754	217,860	9,155	224,497	9,225

cosine similarity between their embedding vectors. When computing these embeddings, the subjects in the sentences are replaced with the same [MASK] tokens so that we only compare the expressed properties. Then, we normalize the logarithmic IDF values into the range [0, 1] to get the distinctiveness scores of clusters.

Specificity. We compute the specificity of an assertion based on the fraction of nouns in it. Concretely, in CANDLE, the specificity of a cluster is computed as the specificity of its summary sentence.

Domain relevance. For each facet, we compute the domain relevance of a cluster by taking the average of the probability scores given to its members by the cultural facet classifier.

Combined score. The final interestingness score for cluster cls is the average of the four feature scores. A higher score means higher interestingness.

Post-filtering. Lastly, to eliminate redundancies and noise, and further improve the final output quality, we employ a few hand-crafted rules:

- At most 500 clusters per subject-facet pair are retained, as further clusters mostly represent redundancies or noise.
- We remove clusters that have no concepts extracted, or that are based on too few distinct sentences (>2/3 same sentences) or web source domains.
- We remove any cluster if either its summary sentence or many of its member sentences match a bad pattern. We compile a set of about 200 regular expression patterns, which were written by a knowledge engineer in one day. For e.g., we reject assertions that contain “the menu”, “the restaurant” (likely advertisements for specific restaurants), or animal and plant breeds named after locations, such as “American bison”, “German Shepherd”, etc.

5 IMPLEMENTATION

Input corpus. In CANDLE, we use the broad web as knowledge source, because of its diversity and coverage, which are important for long-tail subjects. Besides the benefits, the most challenging problem when processing web contents is the tremendous amount of noise, offensive materials, incorrect information etc., hence, choosing a corpus that has been chiefly cleaned is beneficial. We choose the Colossal Clean Crawled Corpus (C4) [27] as our input, a cleaned version of the Common Crawl corpus, created by

Table 2: Processing time and output size of each step in CANDLE for the domain geography>country and facet food.

#	Step	Time	Output/Data size
	Input preprocessing	2 days for NLP	C4 corpus: 8B sentences 196 countries 705 alternate names
1	Subject detection	2 hours	367M subject matches 300M sentences (-96%)
2	Generic assertion filtering	2 hours	13M generic sentences (-96%)
3	Cultural facet classification	4 hours	769K positive sentences (-94%)
4	Assertion clustering	4 hours	42K clusters (-93%)
5	Concept extraction	< 5 minutes	12.4K concepts
6	Cluster ranking and post-filtering	< 5 minutes	8.8K clusters (-80%)
-	Total	~ 12 hours	

applying filters such as deduplication, English-language text detection, removing pages containing source code, offensive language, too little content, etc. We use the C4.EN split, which contains 365M English articles, each with text content and source URL. Before passing it to our system, we preprocessed all C4 documents using spaCy, which took 2 days on our cluster of 6K CPU cores.

Subjects. We collect CCSK for subjects from 3 cultural domains: geography (272 subjects), religions (14 subjects) and occupations (100 subjects). For geography, we split into 4 sub-domains: countries, continents, geopolitical regions (e.g., Middle East, Southeast Asia, etc.) and US states, which were collected from the GeoNames database³, which also provides alias names. We further enriched these aliases with demonyms from Wikipedia⁴.

Facets of culture. We consider 5 facets: *food, drinks, clothing, rituals, and traditions* (for geography/religion) or *behaviors* (for occupation), selected based on an article on facets of culture [23].

Execution and result statistics. After tuning the system’s hyperparameters on small withheld data (see Appx. B), we executed CANDLE on a cluster of 6K CPU cores (AMD EPYC 7702) and 40 GPUs (a mix of NVIDIA RTX 8000, Tesla A100 and A40 GPUs).

Regarding processing time, for the domain country (196 subjects), it took a total of 12 hours to complete the extraction, resulting in 8.4K clusters for the facet food (cf. Table 2). Occupations and religions took 8 and 6 hours each.

We provide statistics of the output in Table 1. In total, the resulting collection has 1.1M CCSK assertions (i.e., base sentences) which form 60K clusters for the given subjects and facets.

6 EVALUATION

We perform the following evaluations:

- (1) A comparison of quality of CANDLE’s output and existing socio-cultural CSK resources: This analysis will show that

our CCSK collection is of significantly higher quality than existing resources (Sec. 6.1), and even outperforms GPT-3-generated assertions (Sec. 6.2).

- (2) Two extrinsic use cases for CCSK: In this evaluation, we perform two downstream applications, question answering (QA) and a “guess the subject” game, showing that using CCSK assertions from CANDLE is beneficial for these tasks, and that CANDLE assertions outperform those generated by GPT-3 (Sec. 6.3).

In Appx. C, we also break down our CCSK collection into domains and facets, analyzing in details the assertion quality for each sub-collection.

6.1 Comparison with other resources

6.1.1 Evaluation metrics. Following previous works [7, 30], we analyze assertion quality along several complementary metrics, annotated by Amazon MTurk (AMT) crowdsourcing.

- (1) **Plausibility (PLA).** This dimension measures whether assertions are considered to be generally true, a CCSK-softened variant of correctness/precision.
- (2) **Commonality (COM).** This dimension measures whether annotators have heard of the assertion before, as a signal for whether assertions cover mainstream or fringe knowledge (akin to salience).
- (3) **Distinctiveness (DIS).** This dimension measures discriminative informativeness of assertions, i.e., whether the assertion differentiates the subject from others.

Each metric is evaluated on a 3-point Likert scale for negation (0), ambiguity (1) and affirmation (2). Distinctiveness (DIS) is only applicable if the answer to the plausibility (PLA) question is either 1 or 2. In case the annotators are not familiar with the assertion, we encourage them to perform a quick search on the web to find out the answers for the PLA and DIS questions. More details on the AMT setup can be found in Appx. D.

6.1.2 Compared resources. We compare CANDLE with 3 prominent CSK resources: Quasimodo [30], Acharya et al. [1], StereoKG [7]. The former covers broad domains including assertions for countries and religions, while the others focus on cultural knowledge. Other popular resources such as ConceptNet [36], GenericsKB [3], Ascent/Ascent++ [21, 22], ATOMIC [32], ASER [47] and TransOMCS [45] do not have their focus on cultural knowledge and contain very little to zero assertions for geography or religion subjects, hence, they are not qualified for this comparison.

We evaluate 2 versions of CANDLE, one where each base assertion is retained independently (CANDLE-base-sent), the other containing only the cluster representatives (CANDLE-cluster-reps).

6.1.3 Setup. For comparability, all resources are compared on 100 random assertions of the same 5 country subjects covered in StereoKG [7] - United States, China, India, Germany and France. We note that among all compared resources, Acharya et al. [1] only contain two subjects (United States and India), so for that resource, we only sample from those. For StereoKG, we use their natural-language assertions. For Quasimodo and Acharya et al., we verbalize their triples using crafted rules. Each assertion is evaluated by 3 MTurk annotators. Additionally, we ask if the annotator would

³<http://www.geonames.org/>

⁴<https://en.wikipedia.org/wiki/Demonym>

Table 3: CANDLE in comparison to other CSK resources. Quality evaluated on assertions of 5 popular countries in StereoKG. Abbrev.: PLA - plausibility, COM - commonality, DIS - distinctiveness, OFF - offensiveness, LEN - average assertion length.

Resource	Construction	Format	Size			Quality [0..2]			OFF (%)	LEN
			196 countries	10 religions	100 occupations	PLA	COM	DIS		
Acharya et al. [1]	Crowdsourcing	Fixed relations	225	0	0	1.32	1.22	0.25	2	102
StereoKG [7]	Text extraction	OpenIE triples	2,181	1,810	0	0.54	0.46	0.21	18	37
Quasimodo [30]	Text extraction	OpenIE triples	22,588	10,628	51,124	0.68	0.65	0.31	13	32
CANDLE-base-sent	Text extraction	Sentences	520,971	226,807	238,057	1.21	0.93	0.76	1	69
CANDLE-cluster-reps	Text extraction	Sentences	28,711	8,823	9,826	1.50	1.15	1.03	1	73

consider the assertion as an inappropriate or offensive material. More details on the annotation task can be found in Appx. D.

6.1.4 Results. A summary of comparison with other resources is shown in Table 3.

Resource size and assertion length. CANDLE outperforms all other resources on the number of base sentences. When turning to clusters, our resource still has significantly more assertions than Acharya et al. (which was constructed manually at small scale) and StereoKG (extracted from Reddit/Twitter questions). Quasimodo has comparable size with CANDLE-cluster-reps for the country and religion domains and has more for the occupation domain.

The OpenIE-based methods, Quasimodo and StereoKG, produce the shortest assertion (32 and 37 characters on average, respectively). The manually-constructed KG (Acharya et al.) has the longest assertions (102 characters). CANDLE, having average assertion lengths (69 and 73), stands between those two approaches.

Assertion quality. In general, CANDLE-cluster-reps considerably outperforms all other baselines on 2 of the 3 metrics (plausibility and distinctiveness). Our resource only comes behind Acharya et al. on the *commonality* metric (1.15 and 1.22 respectively), which is expected because Acharya et al. only cover a few relations about common rituals (e.g., birthday, wedding, funeral) in two countries, USA and India, and their assertions are naturally known by many workers on Amazon MTurk, who are mostly from these 2 countries [31]. Importantly, the resource of Acharya et al. is based on crowdsourcing and only contains a small set of 225 assertions for a few rituals.

CANDLE-cluster-reps even outperforms the manually-constructed KG (Acharya et al.) on the *plausibility* metric. This could be caused by an annotation task design that is geared towards abnormalities, or lack of annotation quality assurance.

CANDLE also has the highest scores on the *distinctiveness* metric, while most of the assertions in other resources were marked as not distinguishing by the annotators.

Between the two versions of CANDLE, the cluster representatives consistently outperform the base sentences on all evaluated metrics. This indicates that still some of the raw sentences in the collection are noisy, on the other hand, the computed cluster representatives are more coherent and generally of better quality.

We also measured the *offensiveness* (OFF) of each resource, i.e., the percentage of assertions that were marked as inappropriate or offensive materials by at least one of the human-annotators. Quasimodo and StereoKG, extracted from raw social media contents, have the highest number of assertions considered offensive (18%

Table 4: Assertion quality - CANDLE vs. GPT-3 - evaluated on assertions of 196 countries.

Method	Quality [0..2]			OFF (%)	LEN
	PLA	COM	DIS		
GPT-3 [5]	1.26	0.80	0.73	1	81
CANDLE	1.25	0.89	0.89	1	75

and 13%). Meanwhile, CANDLE’s judicious filters only miss a small fraction (1% of final assertions).

In summary, our CANDLE CCSK collection has the highest quality by a large margin compared to other resources. Our resource provides assertions of high plausibility and distinctiveness. The clustering and cluster summarization also help to improve the presentation quality of the CCSK.

6.2 Comparison with direct LM extraction

Knowledge extraction directly from pre-trained LMs is recently popular, e.g., the LAMA probe [25] or AutoTOMIC [40]. There are major pragmatic challenges to this approach, in particular, that assertions cannot be contextualized with truly observed surrounding sentences, and that errors cannot be traced back to specific sources. Nonetheless, it is intrinsically interesting to compare assertion quality between extractive and generative approaches. In this section, we compare CANDLE with assertions generated by the state-of-the-art LM, GPT-3 [5].

Generating knowledge with GPT-3. We query the largest GPT-3 model (*davinci-002*) with the following prompt template: “Please write 20 short sentences about notable <facet> in <subject>.” We run each prompt 10 times and set the randomness (temperature) to 0.7, so as to obtain a larger resource. We run the query for 5 facets and 210 subjects (196 countries and 14 religions), resulting in 188,061 unique sentences. Henceforth we call this dataset *GPT-resource*, and reuse it in the extrinsic use cases (Sec. 6.3).

Evaluation metrics and setup. For each resource, we sample 100 assertions for each facet (hence, 500 assertions in total) and perform human evaluation on the 3 metrics - commonality (COM), plausibility (PLA) and distinctiveness (DIS).

Results. The quality comparison between assertions of CANDLE and *GPT-resource* is shown in Table 4. While plausibility scores are the same, and CANDLE performs better in commonality, the difference that stands out is in distinctiveness: GPT-3 performs

Table 5: Example assertions of CANDLE and GPT-resource for subject:China, facet:clothing.

#	CANDLE	GPT-resource
1	The bride usually wears red in a traditional Chinese wedding.	Chinese people also like to wear modern clothes such as jeans and t-shirts.
2	The Chinese wear white at funerals bec. it is associated with mourning in Chinese culture.	Shoes are also very important in Chinese culture.
3	The Chinese wear new clothes for the New Year to symbolize new beginnings.	Chinese people also like to dress their children in very cute clothes.
4	The costumes in Chinese opera are very colorful and important.	In China, you will often see little girls wearing dresses and boys wearing shorts.
5	In ancient China, only the emperor was allowed to wear the color yellow.	In the winter, people in China wear coats and scarves to keep warm.

Table 6: Results of QA using context-augmented LMs.

Facet	#Questions	Precision (%)		
		No cont.	GPT cont.	CANDLE cont.
Food/Drinks	88	92.05	94.32	93.18
Behaviors	125	60.80	57.60	63.20
Rituals	135	87.41	85.93	92.59
Traditions	152	72.37	69.74	79.61
All	500	77.00	75.40	81.40

significantly worse, reconfirming a known problem of language models, evasiveness and over-generality [17]. We illustrate this with anecdotal evidence in Table 5, for subject:China and facet:clothing. None of the listed GPT-3 examples is specific for China.

6.3 Extrinsic evaluation

QA with context-augmented LMs. Augmenting LMs input with additional contexts retrieved from knowledge bases has been a popular approach to question answering (QA) [11, 24], which shows that although LMs store information in billions of parameters, they still lack knowledge to answer knowledge-intensive questions, e.g., “What is the appropriate color to wear at a Hindu funeral?”

In this experiment, we use GPT-3 as QA agent, and compare its performance in 3 settings: (1) when only the questions are given, and when questions and their related contexts retrieved from (2) CANDLE or (3) GPT-resource (cf. Sec. 6.2) are given to the LM. For *questions*, we collect cultural knowledge quizzes from multiple websites, which results in 500 multiple-choice questions, each with 2-5 answer options (only one of them is correct). For *context retrieval*, we use the SentenceBert *all-mpnet-base-v2* model, and for each question, retrieve the two most similar assertions from CANDLE-cluster-reps and GPT-resource. We use the GPT-3 *davinci-002* model, with temperature=0 and max_length=16 (see Appx. E for prompt settings).

We measure the precision of the answers and present the results in Table 6. It can be seen that with CANDLE context, the performance is consistently better than when no context is given on all facets of culture, and better than GPT context on 3 out of 4 facets. This shows that GPT-3, despite its hundred billions of parameters, still lacks socio-cultural knowledge for question answering, and external resources such as CANDLE CCSK can help to alleviate this problem.

“Guess the country” game. The rule of this game is as follows: Given 5 CCSK assertions about a country, a player has to guess the name of the country.

Table 7: Precision (%) for the “guess the country” game.

	Food	Drinks	Clothing	Rituals	Traditions	Avg.
GPT-resource	63.0	30.0	44.0	70.0	84.0	58.2
CANDLE	85.0	74.0	62.0	76.0	80.0	75.4

As *input*, we select a random set of 100 countries, and take assertions from either CANDLE or GPT-resource. The game has 5 rounds, each is associated with a facet of culture. In each round, for each country, we draw the top-5 assertions from each resource (sorted by interestingness in CANDLE or by frequency in GPT-resource). All mentions of the countries in the input sentences are replaced with [. . .], before being revealed to the player.

This is a game that requires a player that possesses a wide range of knowledge across many cultures. Instead of human players, we choose GPT-3 as our player, which has been shown to be excellent at many QA tasks [5] (prompt settings are presented in Appx. E).

We measure the precision of the answers and present the results in Table 7. It can be seen that the player got significantly more correct answers when given assertions from CANDLE than from GPT-resource (i.e., assertions written by the player itself!). This confirms that assertions in CANDLE are more informative.

7 CONCLUSION

We presented CANDLE—an end-to-end methodology for automatically collecting cultural commonsense knowledge (CCSK) from broad web contents at scale. We executed CANDLE on several cultural subjects and facets of culture and produce CCSK of high quality. Our experiments showed the superiority of the resulting CCSK collection over existing resources, which have limited coverage for this kind of knowledge, and also over methods based on prompting LMs. Our work expands CSKG construction into a domain that has been largely ignored so far. Our data and code are accessible at <https://cultural-csk.herokuapp.com/>.

Ethics statement

No personal data was processed and hence no IRB review was conducted. It is in the nature of this research, however, that some outputs reflect prejudices or are even offensive. We have implemented multiple filtering steps to mitigate this, and significantly reduced the percentage of offensive assertions, compared with prior work. Nonetheless, CANDLE represents a research prototype, and outputs should not be used in downstream tasks without further thorough review.

REFERENCES

- [1] Anurag Acharya, Kartik Talamadupula, and Mark A. Finlayson. 2020. An Atlas of Cultural Commonsense for Machine Reasoning. *CoRR* abs/2009.05664 (2020), 9 pages. arXiv:2009.05664 <https://arxiv.org/abs/2009.05664>
- [2] Junia Anacleto, Henry Lieberman, Marie Tsutsumi, Vânia Neris, Aparecido Carvalho, Jose Espinosa, Muriel Godoi, and Silvia Zem-Mascarenhas. 2006. Can Common Sense uncover cultural differences in computer applications?. In *Artificial Intelligence in Theory and Practice*, Max Bramer (Ed.). Springer US, Boston, MA, 1–10. https://doi.org/10.1007/978-0-387-34747-9_1
- [3] Sumithra Bhakthavatsalam, Chloe Anastasiades, and Peter Clark. 2020. GenericsKB: A Knowledge Base of Generic Statements. *CoRR* abs/2005.00660 (2020), 6 pages. arXiv:2005.00660 <https://arxiv.org/abs/2005.00660>
- [4] Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 4762–4779. <https://doi.org/10.18653/v1/P19-1470>
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models Are Few-Shot Learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (Vancouver, BC, Canada) (NIPS'20)*. Curran Associates Inc., Red Hook, NY, USA, Article 159, 25 pages. <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>
- [6] Bhavana Dalvi Mishra, Niket Tandon, and Peter Clark. 2017. Domain-Targeted, High Precision Knowledge Extraction. *Transactions of the Association for Computational Linguistics* 5 (2017), 233–246. https://doi.org/10.1162/tacl_a_00058
- [7] Awantee Deshpande, Dana Ruitter, Marius Mosbach, and Dietrich Klakow. 2022. StereoKG: Data-Driven Knowledge Graph Construction For Cultural Knowledge and Stereotypes. In *Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH)*. Association for Computational Linguistics, Seattle, Washington (Hybrid), 67–78. <https://aclanthology.org/2022.woah-1.7>
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [9] Joseph L. Fleiss and Jacob Cohen. 1973. The Equivalence of Weighted Kappa and the Intraclass Correlation Coefficient as Measures of Reliability. *Educational and Psychological Measurement* 33, 3 (1973), 613–619. <https://doi.org/10.1177/001316447303300309>
- [10] Jonathan Gordon, Benjamin Van Durme, and Lenhart K. Schubert. 2010. Learning from the Web: Extracting General World Knowledge from Noisy Text. In *Proceedings of the 2nd AAI Conference on Collaboratively-Built Knowledge Sources and Artificial Intelligence (AAIWS'10-02)*. AAAI Press, Palo Alto, California, USA, 10–15. <https://www.aaai.org/ocs/index.php/WS/AAIW10/paper/viewFile/2035/2408>
- [11] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: Retrieval-Augmented Language Model Pre-Training. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*. JMLR.org, Online, Article 368, 10 pages. <http://proceedings.mlr.press/v119/guu20a/guu20a.pdf>
- [12] Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia D'amato, Gerard De Melo, Claudio Gutierrez, Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian Neumaier, Axel-Cyrille Ngonga Ngomo, Axel Polleres, Sabbir M. Rashid, Anisa Rula, Lukas Schmelzeisen, Juan Sequeda, Steffen Staab, and Antoine Zimmermann. 2021. Knowledge Graphs. *ACM Comput. Surv.* 54, 4, Article 71 (jul 2021), 37 pages. <https://doi.org/10.1145/3447772>
- [13] Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. (Comet-)Atomic 2020: On Symbolic and Neural Commonsense Knowledge Graphs. In *Thirty-Fifth AAI Conference on Artificial Intelligence (AAAI'21)*. AAAI Press, Palo Alto, California, USA, 6384–6392. <https://ojs.aaai.org/index.php/AAAI/article/view/16792>
- [14] Filip Ilievski, Pedro Szekely, and Bin Zhang. 2021. CSKG: The CommonSense Knowledge Graph. In *The Semantic Web: 18th International Conference, ESWC 2021, Virtual Event, June 6–10, 2021, Proceedings*. Springer-Verlag, Berlin, Heidelberg, 680–696. https://doi.org/10.1007/978-3-030-77385-4_41
- [15] Douglas B. Lenat. 1995. CYC: A Large-Scale Investment in Knowledge Infrastructure. *Commun. ACM* 38, 11 (Nov 1995), 33–38. <https://doi.org/10.1145/219717.219745>
- [16] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 7871–7880. <https://doi.org/10.18653/v1/2020.acl-main.703>
- [17] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A Diversity-Promoting Objective Function for Neural Conversation Models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, San Diego, California, 110–119. <https://doi.org/10.18653/v1/N16-1014>
- [18] Fangyu Liu, Emanuele Bugliarelli, Edoardo Maria Ponti, Siva Reddy, Nigel Collier, and Desmond Elliott. 2021. Visually Grounded Reasoning across Languages and Cultures. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 10467–10485. <https://doi.org/10.18653/v1/2021.emnlp-main.818>
- [19] Hugo Liu and Push Singh. 2004. ConceptNet — A Practical Commonsense Reasoning Tool-Kit. *BT Technology Journal* 22 (2004), 211–226. <https://doi.org/10.1023/B:BTJ.0000047600.45421.6d>
- [20] Tuan-Phong Nguyen and Simon Razniewski. 2022. Materialized Knowledge Bases from Commonsense Transformers. In *Proceedings of the First Workshop on Commonsense Representation and Reasoning (CSRR 2022)*. Association for Computational Linguistics, Dublin, Ireland, 36–42. <https://doi.org/10.18653/v1/2022.csrr-1.5>
- [21] Tuan-Phong Nguyen, Simon Razniewski, Julien Romero, and Gerhard Weikum. 2022. Refined Commonsense Knowledge from Large-Scale Web Contents. *IEEE Transactions on Knowledge and Data Engineering* 0, 0 (2022), 1–16. <https://doi.org/10.1109/TKDE.2022.3206505>
- [22] Tuan-Phong Nguyen, Simon Razniewski, and Gerhard Weikum. 2021. Advanced Semantics for Commonsense Knowledge Extraction. In *Proceedings of the Web Conference 2021 (Ljubljana, Slovenia) (WWW '21)*. Association for Computing Machinery, New York, NY, USA, 2636–2647. <https://doi.org/10.1145/3442381.3449827>
- [23] Outline of culture. 2022. Outline of culture — Wikipedia, The Free Encyclopedia. https://en.wikipedia.org/wiki/Outline_of_culture Online; accessed: 2022-10-08.
- [24] Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2020. How Context Affects Language Models' Factual Predictions. In *Automated Knowledge Base Construction (AKBC'20)*. AKBC.ws, Online, 15 pages. <https://www.akbc.ws/2020/assets/pdfs/025X0zPfn.pdf>
- [25] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language Models as Knowledge Bases?. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 2463–2473. <https://doi.org/10.18653/v1/D19-1250>
- [26] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9. https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf
- [27] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research* 21, 140 (2020), 1–67. <http://jmlr.org/papers/v21/20-074.html>
- [28] Simon Razniewski, Niket Tandon, and Aparna S. Varde. 2021. Information to Wisdom: Commonsense Knowledge Extraction and Compilation. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining (Virtual Event, Israel) (WSDM '21)*. Association for Computing Machinery, New York, NY, USA, 1143–1146. <https://doi.org/10.1145/3437963.3441664>
- [29] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3982–3992. <https://doi.org/10.18653/v1/D19-1410>
- [30] Julien Romero, Simon Razniewski, Koninika Pal, Jeff Z. Pan, Archit Sakhadeo, and Gerhard Weikum. 2019. Commonsense Properties from Query Logs and Question Answering Forums. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (Beijing, China) (CIKM '19)*. Association for Computing Machinery, New York, NY, USA, 1411–1420. <https://doi.org/10.1145/3357384.3357955>
- [31] Joel Ross, Lilly Irani, M. Six Silberman, Andrew Zaldivar, and Bill Tomlinson. 2010. Who Are the Crowdworkers? Shifting Demographics in Mechanical Turk. In *CHI '10 Extended Abstracts on Human Factors in Computing Systems (Atlanta, Georgia, USA) (CHI EA '10)*. Association for Computing Machinery, New York, NY, USA, 2863–2872. <https://doi.org/10.1145/1753846.1753873>
- [32] Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019. ATOMIC: An Atlas of Machine Commonsense for If-Then Reasoning. In *Proceedings of the Thirty-Third AAI Conference on Artificial Intelligence (Honolulu,*

- Hawaii, USA) (AAAI'19). AAAI Press, Palo Alto, California, USA, Article 372, 9 pages. <https://doi.org/10.1609/aaai.v33i01.33013027>
- [33] Vered Shwartz. 2022. Good Night at 4 pm?! Time Expressions in Different Cultures. In *Findings of the Association for Computational Linguistics: ACL 2022*. Association for Computational Linguistics, Dublin, Ireland, 2842–2853. <https://doi.org/10.18653/v1/2022.findings-acl.224>
- [34] Push Singh, Thomas Lin, Erik T. Mueller, Grace Lim, Travell Perkins, and Wan Li Zhu. 2002. Open Mind Common Sense: Knowledge Acquisition from the General Public. In *On the Move to Meaningful Internet Systems 2002: CoopIS, DOA, and ODBASE*, Robert Meersman and Zahir Tari (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 1223–1237. https://doi.org/10.1007/3-540-36124-3_77
- [35] Karen Sparck Jones. 1988. *A Statistical Interpretation of Term Specificity and Its Application in Retrieval*. Taylor Graham Publishing, GBR, 132–142.
- [36] Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. ConceptNet 5.5: An Open Multilingual Graph of General Knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence* (San Francisco, California, USA) (AAAI'17). AAAI Press, Palo Alto, California, USA, 4444–4451. <https://doi.org/10.1609/aaai.v31i1.11164>
- [37] Niket Tandon, Gerard de Melo, Fabian Suchanek, and Gerhard Weikum. 2014. WebChild: Harvesting and Organizing Commonsense Knowledge from the Web. In *Proceedings of the 7th ACM International Conference on Web Search and Data Mining* (New York, New York, USA) (WSDM '14). Association for Computing Machinery, New York, NY, USA, 523–532. <https://doi.org/10.1145/2556195.2556245>
- [38] Joe H. Ward. 1963. Hierarchical Grouping to Optimize an Objective Function. *J. Amer. Statist. Assoc.* 58 (1963), 236–244.
- [39] Gerhard Weikum, Xin Luna Dong, Simon Razniewski, and Fabian M. Suchanek. 2021. Machine Knowledge: Creation and Curation of Comprehensive Knowledge Bases. *Found. Trends Databases* 10, 2–4 (2021), 108–490. <https://doi.org/10.1561/19000000064>
- [40] Peter West, Chandra Bhagavatula, Jack Hessel, Jena Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2022. Symbolic Knowledge Distillation: from General Language Models to Commonsense Models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Seattle, United States, 4602–4625. <https://aclanthology.org/2022.naacl-main.341>
- [41] Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 1112–1122. <https://doi.org/10.18653/v1/N18-1101>
- [42] Da Yin, Hritik Bansal, Masoud Monajatipoor, Liunian Harold Li, and Kai-Wei Chang. 2022. GeoMLAMA: Geo-Diverse Commonsense Probing on Multilingual Pre-Trained Language Models. *CoRR abs/2205.12247* (2022), 16 pages. <https://doi.org/10.48550/arXiv.2205.12247>
- [43] Da Yin, Liunian Harold Li, Ziniu Hu, Nanyun Peng, and Kai-Wei Chang. 2021. Broaden the Vision: Geo-Diverse Visual Commonsense Reasoning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2115–2129. <https://doi.org/10.18653/v1/2021.emnlp-main.162>
- [44] Wenpeng Yin, Jamaal Hay, and Dan Roth. 2019. Benchmarking Zero-shot Text Classification: Datasets, Evaluation and Entailment Approach. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3914–3923. <https://doi.org/10.18653/v1/D19-1404>
- [45] Hongming Zhang, Daniel Khashabi, Yangqiu Song, and Dan Roth. 2020. TransOMCS: From Linguistic Graphs to Commonsense Knowledge. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*. International Joint Conferences on Artificial Intelligence Organization, California, USA, 4004–4010. <https://doi.org/10.24963/ijcai.2020/554>
- [46] Hongming Zhang, Xin Liu, Haojie Pan, Haowen Ke, Jiefu Ou, Tianqing Fang, and Yangqiu Song. 2022. ASER: Towards large-scale commonsense knowledge acquisition via higher-order selectional preference over eventualities. *Artificial Intelligence* 309 (2022), 103740. <https://doi.org/10.1016/j.artint.2022.103740>
- [47] Hongming Zhang, Xin Liu, Haojie Pan, Yangqiu Song, and Cane Wing-Ki Leung. 2020. ASER: A Large-Scale Eventuality Knowledge Graph. In *Proceedings of The Web Conference 2020 (WWW '20)*. Association for Computing Machinery, New York, NY, USA, 201–211. <https://doi.org/10.1145/3366423.3380107>

A GENERIC FILTERING RULES

GenericsKB [3] was built by using a set of 27 hand-crafted lexico-syntactic rules to extract high-quality generic sentences from different text corpora (the ARC corpus, SimpleWikipedia and the

Waterloo crawl of education websites). For example, the lexical rules look for sentences with short length, starting with a capitalized character, having no bad first words (e.g., determiners), ending with a period, having no URL-like snippets, etc. The syntactic rules only accept a sentence if its root is a verb and not the first word, and if there is a noun before the root verb, etc.

CANDLE adopts the GenericsKB rules. However, as GenericsKB only deals with general concepts (e.g., “tree”, “bird”, etc.), some of the rules are not applicable for the cultural subjects that can be named entities. Hence, depending on the subjects and facets, we adaptively modify the rules (by dropping some of them) so that we will not miss out valuable assertions. For instance, for geography, the *has-no-determiners-as-first-word* rule will filter out valuable assertions such as THE CHINESE USE CHOPSTICKS TO EAT THEIR FOOD or THE CURRYWURST IS A TRADITIONAL GERMAN FAST FOOD DISH, and it must be dropped. In another situation, when exploring the “traditions” facet, the *remove-past-tense-verb-roots* rule would be too aggressive as it rejects assertions about past traditions. The rule that rejects sentences with PERSON entities can be used for the geography and occupation subjects, but must not be used for religions, because it will filter out sentences about Buddha or Jesus Christ. Full details are in the published code base⁵.

B HYPERPARAMETER SETTINGS

Based on tuning on small withheld data, we select the following values for hyperparameters and run CANDLE on the C4 dataset with these settings.

For *cultural facet classification* (cf. Sec. 4.3 and Eq. 1), we fix ρ_+ to 0.5 and ρ_- to 0.3.

For *assertion clustering* (cf. Sec. 4.4), we use the SentenceBert model *all-MiniLM-L6-v2* for computing sentence embeddings. For the HAC algorithm, we measure point-wise Euclidean distance of the normalized embeddings. Then, we use the Ward’s linkage [38], with the maximal distance threshold set to 1.5. In the few cases where input sets are larger, we truncate them at 50K sentences per subject-facet pair, since larger inputs only contain further redundancies, that are not worth the cubic effort of clustering. This concerns only 15 out of 386 subjects. For *cluster summarization*, we only consider the 500 most populated clusters for each subject-facet pair with a minimum size of 3 sentences. More details on prompting GPT-3 for cluster summarization can be found in Appx. E.

For *cluster ranking* (cf. Sec. 4.6), we fix θ in Eq. 3 to 0.8.

C INTRINSIC EVALUATION

We break down the CANDLE CCSK collection into domains and facets and evaluate the assertion quality for each of these sub-collections and get more insights into the produced data.

C.1 Per-domain quality

CANDLE contains 3 cultural domains - geography, religion and occupation. For each domain, we sample 100 assertions and perform crowdsourcing evaluation with the 3 metrics - PLA, COM and DIS (cf. SubSec. 6.1.1). We present the evaluation results in Table 8.

⁵https://github.com/cultural-csk/candle/blob/main/candle/pipeline/component_generic_sentence_filter.py

Table 8: Quality of CANDLE assertions for each domain.

Domain	Quality [0..2]			Acceptance rate (%)		
	PLA	COM	DIS	PLA \geq 1	COM \geq 1	DIS \geq 1
Geography	1.52	1.19	1.03	84.00	66.00	61.33
Religion	1.51	1.29	1.22	85.76	74.67	72.00
Occupation	1.59	1.50	1.25	86.67	82.67	73.67
Average	1.54	1.33	1.17	85.44	74.44	69.00

Table 9: Quality of CANDLE assertions for each facet and the domain *geography>country*.

Facet	Quality [0..2]			
	DOM	PLA	COM	DIS
Food	1.42	1.23	0.94	0.97
Drinks	1.51	1.40	1.14	1.19
Clothing	1.49	1.30	1.04	1.07
Rituals	1.45	1.27	1.06	1.20
Traditions	1.42	1.27	1.02	1.11
Average	1.46	1.29	1.04	1.11

Besides the raw scores (0, 1, 2), we also binarize and denote them as acceptance rates, i.e., a score greater than zero means “accept”.

CANDLE achieves a high *plausibility* (PLA) score of 1.54 on average. Performance on this metric is relatively consistent through all domains. Meanwhile, the *commonality* (COM) metric is highest for the occupation domain and lowest for geography domain.

More than 80% of plausible assertions are annotated as *distinctive* (DIS). Religion and occupation assertions perform significantly better than geography’s on this metric. That could be caused by several assertions for geography subjects being correct but too generic (e.g., JAPANESE FOOD IS ENJOYED BY MANY PEOPLE OR GERMAN BEER IS GOOD). On the other hand, religions and occupations are more distinguishing from one another, while countries or geo-regions usually have cultural overlaps.

C.2 Per-facet quality

We select the assertions for the domain country, and for each facet (food, drinks, clothing, traditions, rituals) we sample 100 assertions for crowdsourcing evaluation. Besides commonality (COM), plausibility (PLA) and distinctiveness (DIS), here we introduce one more evaluation metric, domain relevance (DOM), to measure if an assertion talks about the cultural facet of interest. Only when the DOM score is greater than zero, the other metrics will be evaluated. We present the evaluation results in Table 9.

It can be seen that CANDLE maintains good quality on all evaluation metrics. Notably, scores for the DOM metric are consistently high for all facets, suggesting that the enhanced techniques for zero-shot classification work well on our data. Interestingly, the facet *drinks* outperforms all other facets on 3 of the 4 metrics (DOM, PLA and COM), especially for PLA, its score is significantly higher than others. Assertions for *drinks* and *rituals* are also more distinctive than for other facets.

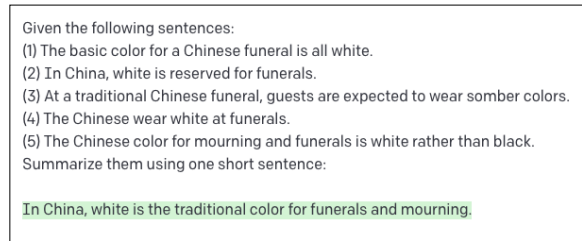


Figure 4: A screenshot of GPT-3 output for cluster summarization.

D DETAILS OF ANNOTATION TASK FOR ASSERTION EVALUATION

The evaluations of assertion quality (Tables 3, 4, 8 and 9) are conducted on Amazon MTurk (AMT). We present CCSK assertions to annotators in the form of natural-language sentences (triples from Quasimodo [30] and Acharya et al. [1] were verbalized using crafted rules). We evaluate each assertion along 3-4 dimensions on a 3-point Likert scale - negation (0), ambiguity (1) and affirmation (2). Each AMT task consists of 5 assertions evaluated by 3 different annotators. Workers are compensated \$0.50 per task. We select Master workers with lifetime’s acceptance rate more than 99%. We obtain fair inter-annotator agreements given by Fleiss’ kappa [9]: 25.0 for DOM, 25.7 for PLA and 25.4 for DIS. This number for COM (13.4) is lower than others because it is an objective question (has the annotator heard of the assertion?).

E GPT-3 PROMPTING

In this work, we use GPT-3 for cluster summarization (Sec. 4.4), generating CCSK for *GPT-resource* (Sec. 6.2), context-augmented QA and “guess the country” game (Sec. 6.3). The prompt templates and settings used for these tasks are presented below.

Cluster summarization. We query the *curie-001* model, with zero temperature and maximum length of 50 tokens. We only take the first generated sentence as output.

Given the following sentences:

- (1) Sentence 1.
- (2) Sentence 2.
- ...
- (n) Sentence n.

Summarize them using one short sentence:

An example prompt is presented in Fig. 4.

Generating CCSK for GPT-resource. We use the largest model (*davinci-002*) and set temperature to 0.7 and maximum length to 512 tokens. For each facet and subject, we run the following prompt template for 10 times: Please write 20 short sentences about notable <facet> in <subject>. We query for 5 facets (food culture, drinking culture, clothing habits, rituals, traditions), and 210 subjects (196 countries and 14 religions). In Table 5, we show some assertions generated using this prompt template for the subject China and the facet “clothing habits”.

Context-augmented QA. We query the *davinci-002* model with zero temperature and maximum length of 16 tokens. Answers are

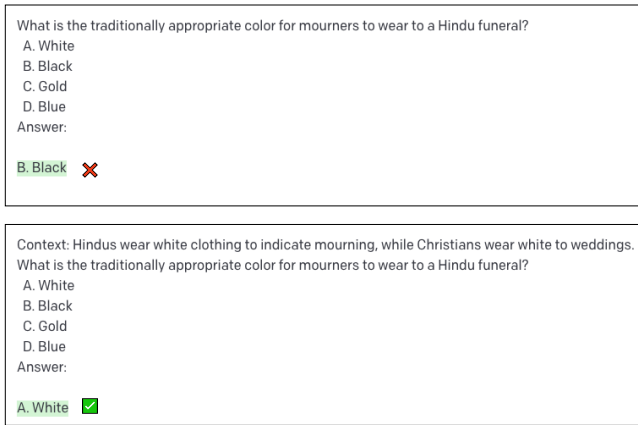


Figure 5: Screenshots of GPT-3 output in the QA task, with-out and with CCSK.

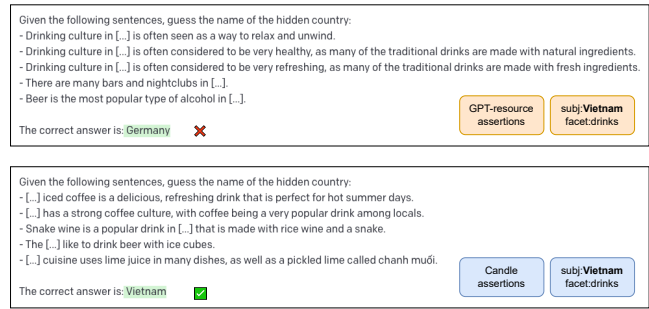


Figure 6: Screenshots of GPT-3 output for the “guess the country” game, with assertions of *GPT-resource* and *CANDLE* for subject: *Vietnam* and facet: *drinks*.

then manually mapped to the respective options. Example prompts are shown in Fig. 5).

“Guess the country” game. We use the *davinci-002* model, with temperature=0 and a max_length=8. Answers given by GPT-3 are checked manually for their correctness. Example prompts can be seen in Fig. 6.