

Information Extraction from PDF Tables with Large Language Models

Anonymous ACL submission

Abstract

Tables, found in PDF documents, contain valuable quantitative information. Unfortunately, extracting this information is difficult due to high variability in the table structure as well as content. We propose statements, a novel data-structure to self-contain quantitative facts and related information. We propose translating tables to statements as a new supervised deep-learning information extraction task. We introduce SemTabNet – a dataset of over 100K annotated tables. Investigating a family of T5-based Statement Extraction Models, our best model predicts statements which are 82% similar to the ground-truth (F1 score of 0.97 for extracting entities). We demonstrate the advantages of representing information as statements by applying our model to over 2700 tables from ESG reports. The homogeneous nature of statements permits data-science analysis on expansive information found in large collections of tables.

1 Introduction

The publishing rate of technical content has increased exponentially (information explosion), in both the academic (Arxiv), legal (USPTO), medical (PubMed) and the commercial domains (annual financial & corporate ESG reports). Many technical documents present their key information in tables. Hence, understanding document tables is important for the field of information extraction.

Large Language Models (LLMs) have been shown to be excellent tools for information extraction, due to their ability to parse, understand, and reason over textual data (OpenAI et al., 2023; Touvron et al., 2023). This, in combination with their ability with zero-shot learning, makes them excellent in information extraction from text (Brown et al., 2020). This approach breaks-down when applying the same techniques on tables (Zhu et al., 2021).

The challenge for understanding tables comes primarily from the high variability in both con-

Multiple Column Headers

Section Headers

Visually Implied Categorical Hierarchy

DESCRIPTION	UNIT	FY20	FY21	FY22
Total community pricing				
Environmental policy	USD	5407,000	5232,919	5751,028
CO2 emissions	tonnes	1,000	1,000	1,000
Percentage of operational sites certified to ISO 14001	%	50%	50%	50%
Percentage of operational sites employees trained on environmental topics	%	100%	100%	100%
DESCRIPTION	UNIT	FY20	FY21	ENVIRONMENTAL GOALS
Total energy consumed (industrial and service sector)	MWh	101,961	111,124	1
Percentage green energy	%	16%	21%	1
Percentage non-green energy	%	84%	79%	1
Total Scope 1 energy	MWh	20,444	25,462	1
Total Scope 2 energy	MWh	76,295	85,662	1
Total energy consumed (industrial and service sector)	MWh	94,246	99,873	1
Non-renewable electricity purchased	MWh	46,223	46,893	1
Renewable electricity purchased	MWh	15,353	19,265	1
Renewable electricity consumed from generated value	MWh	1,989	1,989	1
Consumption of fuel	MWh	27,295	29,893	1
Non-renewable purchased (steam/heating and other purchased energy)	MWh	3,665	4,471	1
Renewable purchased (steam/heating and other purchased energy)	MWh	0	1,242	1
Percentage of total renewable energy	%	16%	21%	1
Energy intensity ratio by floor area	MWh/100 ft²	16.58	17.93	1
Energy intensity ratio per person	MWh/person	1.95	1.95	1
Total value energy will back	\$M	1,000	1,000	1
GHG Emissions				
Scope 1 & 2 (Location-based)	mtCO2e	37,446	39,144	1
Scope 1 & 2 (Market-based)	mtCO2e	22,337	32,912	1
Scope 3 (Location-based)	mtCO2e	6,467	1,111	1
Scope 3 (Market-based)	mtCO2e	29,670	33,793	1
Scope 3 Emissions (Market-based)	mtCO2e	23,203	41,572	1
Total Scope 3 Emissions	mtCO2e	36,137	42,308	1
Category 4: Upstream transportation and distribution	mtCO2e	14,188	19,871	5
Category 5: Business travel	mtCO2e	1,888	2,054	6
Category 6: Employee commuting	mtCO2e	1,787	1,787	7
Category 7: Downstream transportation and distribution	mtCO2e	18,274	18,136	5

Figure 1: Example table from an ESG report with a complicated layout. To extract the information content of a single cell (highlighted in red), the content and relationships (lines drawn in red) to many other cells (highlighted in orange) also needs to be understood.

tent and (spatial) design of document tables. The latter offer a flexible design choice for authors to represent information in a compact format, especially when column and row headers are merged in a hierarchical fashion (see Fig. 1 for an example). This results in a large variability (Kadra et al., 2021; Borisov et al., 2022), with no standardization across domains (e.g. financial reports, corporate ESG reports, scientific papers, patents, books, etc.). Liu et al. (2023) demonstrated that minor perturbations on the structure of a table can seriously undermine the performance of LLMs on downstream tasks. While a lot of progress has been made in table structure recognition, understanding the content of a table is still challenging.

In this paper, we present a general approach for (quantitative) information extraction from tables. First, we propose a new tree-like data-structure, called ‘Statement’, which can combine multiple (named) entities and (n-ary) relations (Fig. 2). It allows us to represent information in a homogeneous domain agnostic fashion. Due to their

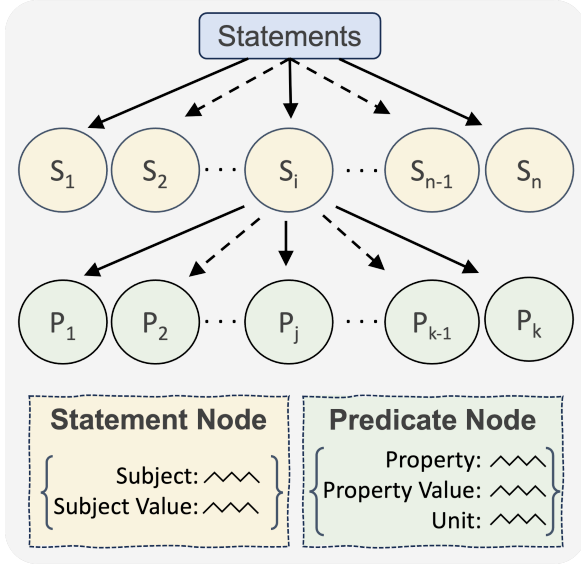


Figure 2: The knowledge model of Statements represented as a tree. From the root node, individual statements emerge as branches. Associated with each individual statement node are the leaf predicate nodes.

knowledge model, statements solve the problem of the variability in table-structure. The nodes of a statement tree can contain content from different subjects, allowing for a general information extraction approach to tables from various domains. Statements can represent information with arbitrary conditions accurately, a must when dealing with complex table layouts containing several multi-column/row headers. With the introduction of statements, the information extraction problem from tables becomes a *translation problem* which we call ‘*statement extraction*’ – translating the original table into a set of statements.

We begin, in Sect. 2 discussing related works. In Sect. 3 we explain the concept of ‘Statements’ and present the SemTabNet dataset used for training our models in Sect. 4. In sect. 5, we discuss the various experiments we performed and their results. We end the paper with an application of our model on ESG reports. Environment, Social, and Governance (ESG) reports which are published by organizations for disclosing their status, and performance on ESG topics. These reports are notoriously hard to parse due to a lack of standardization (Mishra et al., 2023). ESG reports, to this day, are manually analyzed by consultancy firms and professional organisations (Henisz et al., 2019). With our proposed statement extraction, this process can now be fully automated.

2 Related works

LLMs have been widely adopted for information extraction (Xu et al., 2023). Using pre-trained language models, Wang et al. (2022b) perform information extraction in two steps: argument extraction and predicate extraction. Based on this, they introduced a text-based open information extraction benchmark. Wang et al. (2021) presented DeepEx for extracting structured triplets from text based data. Wang et al. (2022a) demonstrate that pre-training models on task-agnostic corpus lead to performance improvement on tasks like information extraction, entity recognition, etc. However, these approaches are limited to textual data.

The application of deep learning to tables has increased due to the availability of large datasets like PubTables-1M (Smock et al., 2021), PubTabNet (Zhong et al., 2020), FinTabNet (Zheng et al., 2020), TabRecSet (Yang et al., 2023), SynthTabNet (Nassar et al., 2022). These datasets focus only on table detection (identifying tables from document images), table structure recognition (parsing table structure) and cell structure recognition (classifying cells as header or data). Additionally, most tables in these datasets are structurally simple, missing out on the complexities of tables encountered in the wild.

A major drawbacks of present approaches is that the semantic meaning of cell content is ignored. This limits the models trained on these datasets. Despite the availability of several attention-based models dedicated to tabular data (TabNet (Arik and Pfister, 2021), TabTransformer (Huang et al., 2020), TableFormer (Nassar et al., 2022), TableFormer (Yang et al., 2022)), Grinsztajn et al. (2022) showed that classic machine learning still performs better than deep neural networks on tabular data.

3 Definition of Statements

The statements data structure aims to homogenize the data representation of information coming from complex, irregular, heterogeneous document tables. At its core, the statements data structure is a tree structure (fig. 2). From the root of the tree, we have ‘subject’-nodes, which contain information regarding the ‘subject’ and the ‘subject-value’ keys. From each subject-node, there are one or more predicate nodes, which define the ‘property’, ‘property-value’, and ‘unit’ keys. Each predicate node carries an atomic piece of quantitative information.

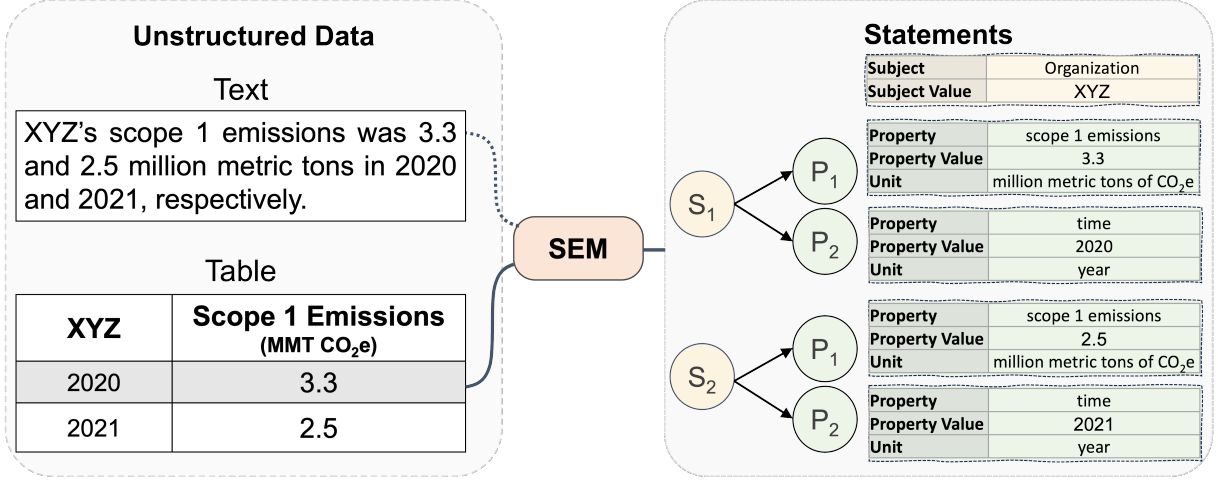


Figure 3: A diagram explaining the framework introduced in this paper. We fine-tune LLMs on the task of ‘Statement Extraction’ leading to a family of “Statement Extraction Models” (SEM). Quantitative facts are extracted from heterogenous unstructured data (only tables in this paper) and stored as Statements.

The statement knowledge model can be applied to both text and tables. In Fig. 3, we show the same statements structure which could be obtained from a text or a corresponding table. As such, the statements structure is not bound only to tables, however, it shows its usefulness particularly when normalising information from heterogeneous tables.

Beyond the uniform layout, statements provide a natural way to quantify how much information any source provides. Simply counting the number of nodes in the tree, provides an estimate on the information richness of given source. A statement is complete when it contains all predicates needed to completely specify objective knowledge pertaining to a subject, i.e. all co-dependent predicates.

The tree structure of statements allows us to quantify, with a single number, the transformation of information from a table. This is accomplished by computing the Tree Editing Distance (Pawlik and Augsten, 2016; Schwarz et al., 2017) between predicted and ground-truth statements. TED is defined as the minimum-cost sequence of node operations that transform one tree into another. Two trees are identical if their TED is 0 and maximally distinct if their normalized TED is 1. Like the Levenshtein distances on strings (Levenshtein, 1966), TED involves three kinds of operations: node insertions, deletions, and renaming. The cost of each operation can be freely defined, which makes this metric both flexible and powerful.

For comparing two statement trees, we setup

strict costs for each edit operation. The predictions are maximally punished for any structural deviation from the ground truth, i.e. deletion and insertion each have a cost of 1. For renaming, we only allow two nodes to be renamed if they are of the same type. If both nodes’ value attribute is of type string, then we calculate a normalized Levenshtein edit distance between the two strings. If both nodes’ value attribute is of numerical type, then the two values are directly compared. In this case, the cost is 0 if the two values are the same, and 1 in all other cases. If the value attribute of both the ground truth and the prediction node is empty, then the cost operation is 0. Normalized TED (\bar{t}) is the ratio of the tree edit distance to the number of edits between two trees. Using the normalized TED, a normalized Tree Similarity score can be computed as $t_s = 1 - \bar{t}$.

It is also instructive to look at the edit types which converted the predicted statements into ground-truth statements. For this, we measure the ratio of edit type to the total number of edits. The ratio of insertions/deletions carries information about the structural similarity. If two trees are structurally similar, the edits are dominated by renaming. While tree-based metrics are sensitive to both entity and relationship extraction, we also evaluate entity extraction. For this, we collect all entities from a statement and count true positives when an entity is found in both model prediction and ground truth, and similarly for true negatives and false positives. Based on these, we measure the standard accuracy, recall and F1 measures.

4 SemTabNet: Statements Data

We used the Deep Search toolkit¹ to collect over 10K ESG reports from over 2000 corporations. Deep Search crawled these PDF reports, converted them into machine readable format, and provided this data along with the metadata of each report in json format.

We compiled a list of important keywords which capture many important concepts in ESG reports (see appendix A). Next, we select only those tables which have some relevance with the keywords. For this we used the following conditions: the ROUGE-L precision (longest common sub-sequence) score between raw data and keywords must be greater than 0.75 and there must be quantitative information in the table.

We need a strategy for understanding the content of a table and extracting statements from it. After manually observing hundreds of table, we decided a two step approach to prepare our ground-truth data. First, we classify all the cells in a table based on the semantic meaning of their content into 16 categories which helps us in constructing statements. For each table, this step creates a ‘labels-table’ with the same shape and structure as the original, but the cells of this labels-table only contain category labels (see fig. 4). Secondly, we create a program which reads both the labels-table and the original table and extracts statements in a rule-based approach. The algorithm is described in appendix B. The 16 labels are:

- Property, Property Value
- Sub-property
- Subject, Subject Value
- Unit, Unit Value
- Time, Time Value
- Key, Key Value
- Header 1, Header 2, Header 3
- Empty, Rubbish

During annotation, all cells of a table are mapped to one of the above labels. For cells which contain information pertaining to more than one label, we pick the labels which is higher in our ordered list of labels. So a cell with content “Revenue (US\$)”, is labelled as property. The ‘property’ and ‘sub-property’ cells always have associated ‘property value’ cell(s). The ‘header’ cells never have an associated value and often divide the table into smaller sections. Empty cells are labelled ‘empty’.

¹Available via: <https://ds4sd.github.io>.

Table 1: Counts of data in SemTabNet³.

TASK	TRAIN	TEST	VAL
SE DIRECT	103,455	11,682	5,445
SE INDIRECT 1D	72,580	8,489	3,821
SE INDIRECT 2D	93,153	22,839	4,903

When a table contain unnecessary parts due to faulty table recovery or non-quantitative information. We label such cells as ‘rubbish’. When a property/property value pair carries supplementary information, those cells are annotated as ‘key’/‘key values’.

Additionally, we observed that most tables can be reasonably classified into three baskets: simple, complex, and qualitative. There are simple tables whose structure cannot be further subdivided into any smaller table. There are complex tables whose structure can be further divided into multiple smaller tables. Finally, there are qualitative tables (like table of contents) which contain little valuable information for our endeavour. We collected about 2,800 tables and found ~ 20% were simple, ~ 20% complex, and ~ 60% were qualitative. We discarded all qualitative tables from any further analysis. To ensure that our data is not biased towards either simple or complex tables, we manually annotated all the cells of 569 simple tables and 538 complex tables. In total, we annotated 1,107 tables (84,890 individual cells) giving rise to 42,982 statements.

We further augmented the annotated tables to create a large training data. We shuffle the rows and columns of tables corresponding to property-values to create new augmented tables, while keeping their contents the same. While this is straightforward for simple tables, special care was taken for complex tables such that only rows/columns which belonged together within a category were shuffled. The maximum number of augmented tables emerging from the shuffling operations was limited to 130, leading to over 120K tables. To promote further research and development, we open source this large dataset of semantic cell annotations as SemTabNet². Table 1 shows the counts of (in)-direct statement extraction in SemTabNet.

²The data can be found here.[LINK](#)

³The counts differ slightly due to the manner in which the final data was harmonized. The SE Indirect 1D data consists of the 84 890 original cells annotated from 1 107 tables. The test/train split of tables for SE Indirect 1D was prepared by stratifying across all cell labels. This split was augmented (as described in text) to prepare data for SE Indirect 2D. The

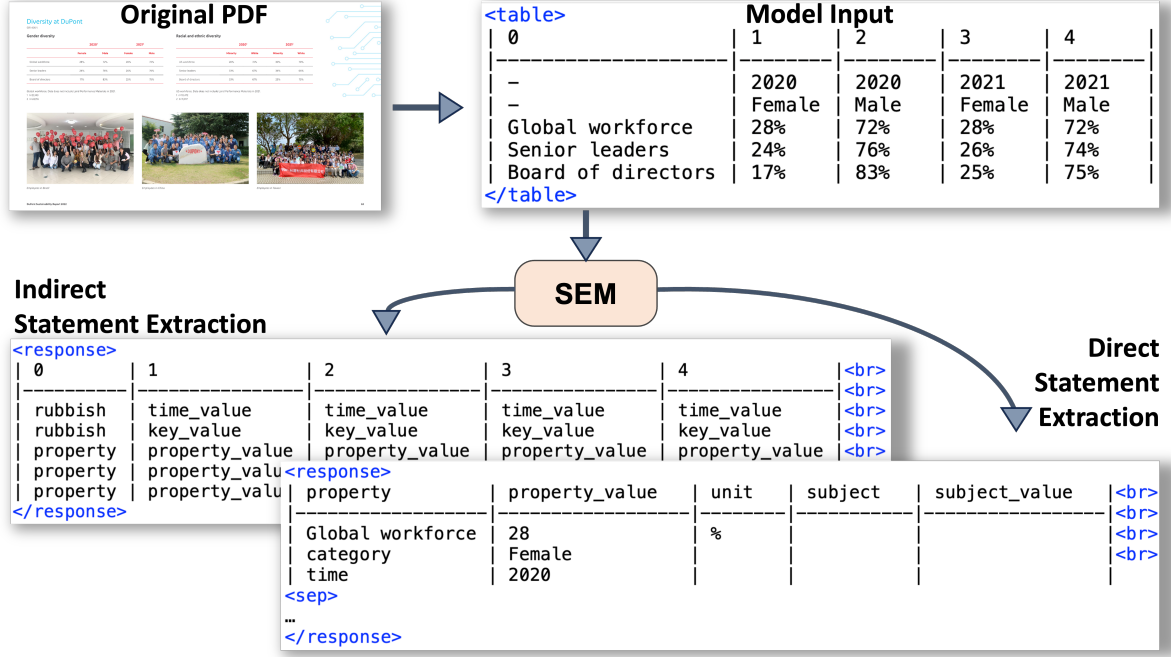


Figure 4: Input and output for the task of “Statement Extraction”. *Top Left*: Page from an ESG report containing tables. *Top Right*: One of the table, from the same page, prepared as markdown for model input. *Bottom Left*: Model output for the task of indirect statement extraction. *Bottom Right*: Model output for the task of direct statement extraction.

5 Experiments & Results

Fig 4 presents Statement Extraction as a supervised deep learning task. Due to the nature of how tables are annotated (see section 4), it is possible to train models for statement extraction statements both directly and indirectly. We consider the following three experiments: (1) *SE Direct*: the model is presented with an input table as markdown in a prompt. The model generates the tabular representation of the resulting statements as markdown. (2) *SE Indirect 1D*: In this experiment, the model input is the individual table cell contents. For a table with n cells, we predict n labels sequentially (hence, 1D) and then use this information to construct statements. Individual cell labels predicted by the model are stitched together to form the labels table, which is then used to construct the predicted statement by using our rule-based algorithm. (3) *SE Indirect 2D*: As opposed to SE Indirect 1D, in this experiment, we predict the cell labels of all cells in a table simultaneously. The entire table, as markdown, is input to the model (hence 2D) and the model generates the labels table, as markdown. Using the rule-based algorithm, the predicted labels

test/train split and augmentation for SE Direct was done independently.

table is converted into predicted statements.

We use six special tokens, which allow us to control and parse model output.

- Input table start token: `<table>`
- Input table stop token: `</table>`
- Output start token: `<response>`
- Output stop token: `</response>`
- Newline token: `
`
- Separate list item token: `<sep>`

This allows us to parse the predicted statements from a LLM. Once successfully parsed, the output statements can be trivially converted from one representation to another. This is crucial because we compare model predicted statements with ground truth by converting statements into a tree structure. These tokens are added to the tokenizer vocabulary before fine-tuning any model.

Since the nature of these tasks naturally fits the paradigm of sequence-to-sequence models, we fine-tune T5 models (Raffel et al., 2020). T5 models are encoder-decoder transformer architecture models which are suitable for many sequence-to-sequence tasks. In our experiments, we train T5 variants (Small, Base, Large, and 3B) to create a family of Statement Extraction Models (SEM).

Table 2: Results of comparing model predicted statements with ground truth data (bold indicates the best in each experiment). For all reported values, the 99% confidence interval, assuming a Gaussian distribution, is $\sim 0.1\%$. The standard error of the mean in all cases is below 0.005%.

Statement Extraction		Context Length	Ratio Tree Edits [%]			Average [%]	
Task	Model		Insert	Delete	Rename	F1	t_s
Indirect 1D	SEM-T5-small	512	98.13	00.00	1.87	62.32	00.86
	SEM-T5-base	512	83.95	01.68	14.37	83.46	09.21
	SEM-T5-large	512	34.68	12.03	53.30	94.67	55.68
	SEM-T5-3b	512	36.70	23.24	40.05	90.49	22.24
Indirect 2D	SEM-T5-small	512	17.34	13.36	69.30	97.06	75.15
	SEM-T5-base	512	15.53	21.60	62.86	96.85	73.87
	SEM-T5-large	512	09.58	22.80	67.62	97.55	80.83
	SEM-T5-3b	512	08.00	28.40	63.59	97.38	81.76
	SEM-T5-small	1024	18.53	18.71	62.75	95.85	68.45
	SEM-T5-base	1024	17.80	16.04	66.16	96.15	69.27
	SEM-T5-large	1024	08.20	17.00	74.79	97.53	79.89
Direct	SEM-T5-small	512	98.14	00.04	01.82	60.65	00.62
	SEM-T5-base	512	97.86	00.06	02.09	68.62	04.46
	SEM-T5-large	512	98.18	00.02	01.80	67.41	04.23
	SEM-T5-3b	512	97.98	0.01	02.01	70.06	03.47
	SEM-T5-small	1024	92.93	00.14	06.93	70.35	02.98
	SEM-T5-base	1024	88.42	00.22	11.35	76.99	11.11
	SEM-T5-large	1024	89.34	00.21	10.45	76.59	06.06

In our training data for tables, the input token count is less than 512 for 50% of the data, and it is less than 1024 for 90% of the data. Thus, except where mentioned, we train T5 models (small, base, large) with context windows of 512 and 1024, and T5-3b with context window of 512. All models are fine-tuned in a distributed data parallel (DDP) manner simultaneously across 4 GPU devices (Nvidia A100-40GB for T5-Small, T5-Base, T5-Large and NVIDIA A100-80GB for T5-3B). Additionally, the largest possible batch size was used for all models. The batch size is impacted by factors like model size, GPU memory, and context window. In turn it affects the number of epochs we can fine-tune in a reasonable time.

For all tasks, we stop the fine-tuning process either after 500,000 steps or after 7 days. We use the AdamW optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. All models are trained with a maximum learning rate of 5×10^{-4} . There is a warm-up phase of 1000 steps in which the learning rate increases linearly from 10^{-10} to 5×10^{-4} . After another 1000 steps, the learning rate is exponentially decayed until it reaches its lowest value of 10^{-6} , where it remains until the end of the training.

Table 2 presents the key results of our exper-

iments. For each table, we evaluate the statements predicted by the model (directly or indirectly) against the ground truth statements. For each task and each model therein, we present the averaged tree similarity score (t_s) (measuring entity & relationship extraction) and the averaged F1 score (measuring entity extraction). Also present are the averaged ratios of tree edit types, which helps us understand t_s . For all reported values, assuming a normal distribution, the standard error of the mean is below 5×10^{-5} and the 99% confidence interval for all values is about $\sim 0.1\%$.

Statement Extraction Indirect 1D: All models trained on this task have context window of 512. Their performance tends to scale with model size. Although the F1 score for entity extraction shows promising value, the tree similarity score for all models is poor. This implies that while these models can learn to extract entities, relationship extraction is difficult for these models. The ratio of tree edits helps us understand these scores. For SEM-T5-small, the ratio of insertion is $\approx 98\%$ which means that the predicted statements does not have enough nodes. As the model size increases, the insertion ratio decreases, and the deletion and renaming ratio increases. Thus, increasing the model

size improves the structural similarity of the predicted tree but the overall performance remains unacceptable.

Statement Extraction Indirect 2D: All models trained on this task perform well on entity extraction with average F1 scores of over 95%. The highest performing model is the SEM-T5-3b (512) with an average tree similarity score of 81.76% and an F1 score of 97.38%. The ratio of insertion edits for this model is the lowest amongst all models across all three tasks. Since the most type of edits required for all model’s predictions in this task are renaming, it implies that the predicted tree has similar structure to the ground truth.

Models with large context window have similar performance on entity extraction but do not perform well on entity and relationship extraction. Since model training cost is quadratic to the sequence length, and we allocate equal training resources to all models, this explains why models with 1024 context window do not show improved performance. We believe that with further training, these models may show better performance than reported here.

Statement Extraction Direct: Based on tree similarity score, most models show poor performance in direct SE. The best performing model, given our training constraints, is SEM-T5-base with a context window of 1024. It gets an average F1 score of 76.99% and an average tree similarity score of only 11%. To understand, why these models performs so poorly on direct SE, we look at the ratio of tree edits.

We note that the ratio of deletions for all models in this task is close to 0. On the other hand, the ratio of insertions for all models is high (from 88% to 98%). This suggests that the statement trees produced by these models is missing vast number of nodes compared to the ground truth. In fact, perusing the model output shows that while the output is of high quality, it contains significantly less nodes than ground truth statements.

Discussion: SE Indirect 1D shows good performance on entity extraction, but performs poorly for both entity and relationship extraction. In this task, the model only sees the content of one cell at a time which makes it easy to extract entities. However, this does not allow the model to develop a strong capability to learn tabular relationships. On the other hand, SE Direct, gives poor performance on both entity extraction and relationship extraction. Direct SE expects the models to unravel a

dense table into statements, for which they must produce many output tokens. For example, the average number of output tokens in the test data for SE direct is 5773 ± 51 , which is significantly larger than the number of tokens for SE indirect 2D (346 ± 1). Thus, direct SE is a very challenging task and might require different strategies to be executed successfully.

SE Indirect 2D, avoids the disadvantages of both the tasks. In this case, the model sees the entire input table (has the chance to learn tabular relationships) and is only tasked with producing a labels table (can finish generation in a reasonable number of tokens). Our experiments clearly demonstrate that statement extraction via the Indirect 2D approach gives better results. This is an unexpected finding of our study, and we hope it motivates other researchers to improve zero-shot statement extraction capability.

6 Application to ESG results

Due to their homogeneous structure, statements enable large-scale exploratory analysis and data science. To demonstrate the advantage of statements over traditional tabular data science, we applied SEM-T5-large (512) over 2700 tables published in over 1000 ESG reports in 2022 using the SE Indirect 2D methodology. This lead to 14,766 statements containing over 100k predicates. This dataset containing ESG related KPIs is invaluable to researchers, policy-makers, and analysts.

We filter this large dataset to contain only those predicates with quantitative property values. This subset contains 47901 predicates from 601 corporate ESG reports. We search the properties in this dataset for some keywords representative of ESG KPIs. Fig. 5 (top) shows the distribution of the number of predicates and the number of distinct organizations which matched our simple keyword search. For example, using ‘emission’ as a keyword, we obtain over 4000 hits with results coming from over 300 distinct corporations. This demonstrates that statements allowed us to pull data from multiple sources and homogenize it for down-stream consumption.

Some of the common properties in this subset are also shown along with their frequency in fig. 5. This shows the breadth and diversity of the nature of information we pulled out from a large corpus of documents. Many of these properties are important ESG KPIs. The ability to extract homogeneously

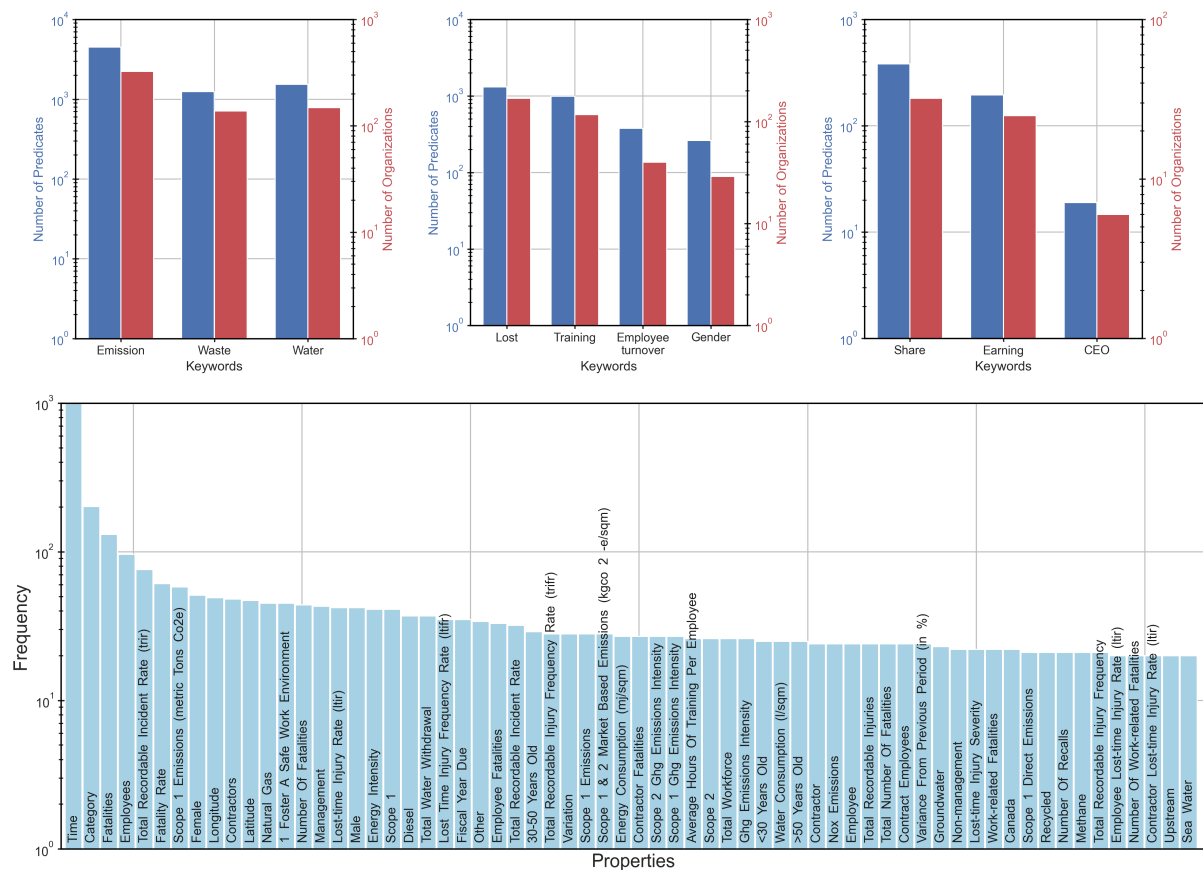


Figure 5: Exploratory analysis of statements from over 2700 Tables published in ESG reports in 2022. *Top:* We searched about 50,000 predicates using keywords (shown on the x-axis) related to environment (left), social (middle), and governance (right). The plot shows the distribution of predicates and the number of organizations from this search. *Bottom:* Frequency chart demonstrating some of the common properties found in our data. This properties are important KPIs in the ESG domain and represent an invaluable data to all stakeholders.

information from a large collection of PDF reports demonstrates the advantage of the statement extraction framework presented in this paper.

7 Conclusion & Future Works

We have presented a novel approach to map complex, irregular, and heterogeneous information to a uniform structure, Statements. We showed how information extraction can be seen as a supervised deep-learning translation task which we called Statement Extraction. We advance the field of table understanding by open-sourcing SemTabNet. SemTabNet consists of 100K tables wherein all cells are annotated reflecting their semantic content. To the best of our knowledge, this is the first work which focuses on the semantic meaning of tabular data.

Investigating three variations of the statement extraction task, we found that using a model to generate table annotations and then construct

statements produces best results. This approach has the advantage, that it produces hallucination-free homogeneous structured data. Statements are an advantageous vehicle for quantitative factual information. They enable down-stream tasks like data science over a large collection of documents. We extracted over 100K facts (predicates) from only 1000 ESG reports.

This work can be easily extended to include domains other than ESG. It can also be extended towards multi-modality by including text data. We leave for future exploration, the use of statements in downstream tasks like QA or document summarization.

8 Limitations

Although, the ideas and the techniques we describe in this paper are domain agnostic, we limit the scope of this paper to the domain of corporate Environment, Social, and Governance (ESG) reports.

This choice is motivated by two observations. First, corporations report valuable quantitative data regarding their efforts to improve their carbon emissions, working conditions, and company culture in ESG reports. These reports contain valuable information regarding the environmental impact of businesses, and the urgency of climate change motivates us to target this domain. Secondly, there is a large variety and diversity of tabular representations used in these reports. Despite efforts to standardize these reports, this diversity makes the task of extracting information from these documents extremely challenging, motivating our choice.

The scope of this work is limited to declarative, explicit knowledge only. All other kinds of knowledge such as cultural, implicit, conceptual, tacit, procedural, conditional, etc. are ignored. We focus on information which one colloquially refers to as ‘hard facts’. Additionally, we limit the scope of this work to quantitative statements i.e. statements whose property values are numerical quantities. We implement this restriction in the notion that we avoid qualitative statements i.e. statements which are not quantitative.

Our model training strategy was biased against large models. We trained all models for either 500K steps or 7 days using the largest possible batch size. This means smaller models learn more frequently (more epochs) than larger models. However, we do not believe this severely impacted the outcome of our experiments. Our resources were enough to recover well-known trends: improved model performance with model size and context-length.

References

Sercan Ö Arik and Tomas Pfister. 2021. [TabNet: Attentive Interpretable Tabular Learning](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(8):6679–6687. Number: 8.

Vadim Borisov, Tobias Leemann, Kathrin Sessler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. 2022. [Deep Neural Networks and Tabular Data: A Survey](#). *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, Chris 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 *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Leo Grinsztajn, Edouard Oyallon, and Gael Varoquaux. 2022. [Why do tree-based models still outperform deep learning on typical tabular data?](#)

Witold Henisz, Tim Koller, and Robin Nuttall. 2019. [Five ways that ESG creates value](#). *McKinsey Quarterly*.

Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. 2020. [TabTransformer: Tabular Data Modeling Using Contextual Embeddings](#). ArXiv:2012.06678 [cs].

Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. 2021. [Well-tuned Simple Nets Excel on Tabular Datasets](#).

V. I. Levenshtein. 1966. Binary Codes Capable of Correcting Deletions, Insertions and Reversals. *Soviet Physics Doklady*, 10:707.

Tianyang Liu, Fei Wang, and Muhao Chen. 2023. [Rethinking Tabular Data Understanding with Large Language Models](#). ArXiv:2312.16702 [cs].

Lokesh Mishra, Cesar Berrospi, Kasper Dinkla, Diego Antognini, Francesco Fusco, Benedikt Bothur, Maksym Lysak, Nikolaos Livathinos, Ahmed Nassar, Panagiotis Vagenas, Lucas Morin, Christoph Auer, Michele Dolfi, and Peter Staar. 2023. [ESG Accountability Made Easy: DocQA at Your Service](#). ArXiv:2311.18481 [cs].

Ahmed Nassar, Nikolaos Livathinos, Maksym Lysak, and Peter Staar. 2022. [TableFormer: Table Structure Understanding with Transformers](#). In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4604–4613, New Orleans, LA, USA. IEEE.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madeleine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,

647	Simón Posada Fishman, Juston Forte, Isabella Ful-	vin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang	711
648	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	Zhuang, William Zhuk, and Barret Zoph. 2023. GPT-	712
649	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	4 Technical Report . ArXiv:2303.08774 [cs].	713
650	Lopes, Jonathan Gordon, Morgan Grafstein, Scott		
651	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	Mateusz Pawlik and Nikolaus Augsten. 2016. Tree edit	714
652	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	distance: Robust and memory-efficient . <i>Information</i>	715
653	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	<i>Systems</i> , 56:157–173.	716
654	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,		
655	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	717
656	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	718
657	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	Wei Li, and Peter J. Liu. 2020. Exploring the Lim-	719
658	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	its of Transfer Learning with a Unified Text-to-Text	720
659	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	Transformer . <i>Journal of Machine Learning Research</i> ,	721
660	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	21(140):1–67.	722
661	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,		
662	Christina Kim, Yongjik Kim, Hendrik Kirchner,	Stefan Schwarz, Mateusz Pawlik, and Nikolaus Aug-	723
663	Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz	sten. 2017. A New Perspective on the Tree Edit Dis-	724
664	Kondraciuk, Andrew Kondrich, Aris Konstantini-	tance . In <i>Similarity Search and Applications</i> , Lecture	725
665	dis, Kyle Kopic, Gretchen Krueger, Vishal Kuo,	Notes in Computer Science, pages 156–170, Cham.	726
666	Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike,	Springer International Publishing.	727
667	Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim,		
668	Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa	Brandon Smock, Rohith Pesala, and Robin Abra-	728
669	Lopez, Ryan Lowe, Patricia Lue, Anna Makanju,	ham. 2021. PubTables-1M: Towards comprehen-	729
670	Kim Malfacini, Sam Manning, Todor Markov, Yaniv	sive table extraction from unstructured documents .	730
671	Markovski, Bianca Martin, Katie Mayer, Andrew	ArXiv:2110.00061 [cs].	731
672	Mayne, Bob McGrew, Scott Mayer McKinney,		
673	Christine McLeavey, Paul McMillan, Jake McNeil,	Hugo Touvron, Louis Martin, and Kevin Stone. 2023.	732
674	David Medina, Aalok Mehta, Jacob Menick, Luke	Llama 2: Open Foundation and Fine-Tuned Chat	733
675	Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie	Models .	734
676	Monaco, Evan Morikawa, Daniel Mossing, Tong Mu,		
677	Mira Murati, Oleg Murk, David Mély, Ashvin Nair,	Chenguang Wang, Xiao Liu, Zui Chen, Haoyun Hong,	735
678	Reiichiro Nakano, Rajeev Nayak, Arvind Neelakan-	Jie Tang, and Dawn Song. 2021. Zero-Shot Infor-	736
679	tan, Richard Ngo, Hyeonwoo Noh, Long Ouyang,	mation Extraction as a Unified Text-to-Triple Trans-	737
680	Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe	lation . In <i>Proceedings of the 2021 Conference on</i>	738
681	Palermo, Ashley Pantuliano, Giambattista Paras-	<i>Empirical Methods in Natural Language Processing</i> ,	739
682	candolo, Joel Parish, Emy Parparita, Alex Passos,	pages 1225–1238, Online and Punta Cana, Domini-	740
683	Mikhail Pavlov, Andrew Peng, Adam Perelman, Fil-	can Republic. Association for Computational Lin-	741
684	ipe de Avila Belbute Peres, Michael Petrov, Henrique	guistics.	742
685	Ponde de Oliveira Pinto, Michael, Pokorny, Michelle		
686	Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power,	Chenguang Wang, Xiao Liu, Zui Chen, Haoyun Hong,	743
687	Boris Power, Elizabeth Proehl, Raul Puri, Alec	Jie Tang, and Dawn Song. 2022a. DeepStruct: Pre-	744
688	Radford, Jack Rae, Aditya Ramesh, Cameron Ray-	training of Language Models for Structure Prediction .	745
689	mond, Francis Real, Kendra Rimbach, Carl Ross,	In <i>Findings of the Association for Computational Lin-</i>	746
690	Bob Rotsted, Henri Roussez, Nick Ryder, Mario	<i>guistics: ACL 2022</i> , pages 803–823, Dublin, Ireland.	747
691	Saltarelli, Ted Sanders, Shibani Santurkar, Girish	Association for Computational Linguistics.	748
692	Sastry, Heather Schmidt, David Schnurr, John Schul-		
693	man, Daniel Selsam, Kyla Sheppard, Toki Sherbakov,	Chenguang Wang, Xiao Liu, and Dawn Song. 2022b.	749
694	Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon	IELM: An Open Information Extraction Benchmark	750
695	Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin,	for Pre-Trained Language Models . In <i>Proceedings of</i>	751
696	Katarina Slama, Ian Sohl, Benjamin Sokolowsky,	<i>the 2022 Conference on Empirical Methods in Nat-</i>	752
697	Yang Song, Natalie Staudacher, Felipe Petroski Such,	<i>ural Language Processing</i> , pages 8417–8437, Abu	753
698	Natalie Summers, Ilya Sutskever, Jie Tang, Niko-	Dhabi, United Arab Emirates. Association for Com-	754
699	las Tezak, Madeleine Thompson, Phil Tillet, Amin	putational Linguistics.	755
700	Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick		
701	Turley, Jerry Tworek, Juan Felipe Cerón Uribe, An-	Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong	756
702	drea Vallone, Arun Vijayvergiya, Chelsea Voss, Car-	Xu, Xiangyu Zhao, Xian Wu, Yefeng Zheng, and	757
703	roll Wainwright, Justin Jay Wang, Alvin Wang, Ben	Enhong Chen. 2023. Large language models for	758
704	Wang, Jonathan Ward, Jason Wei, C. J. Weinmann,	generative information extraction: A survey .	759
705	Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian		
706	Weng, Matt Wiethoff, Dave Willner, Clemens Win-	Fan Yang, Lei Hu, Xinwu Liu, Shuangping Huang, and	760
707	ter, Samuel Wolrich, Hannah Wong, Lauren Work-	Zhenghui Gu. 2023. A large-scale dataset for end-	761
708	man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao,	to-end table recognition in the wild . <i>Scientific Data</i> ,	762
709	Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Woj-	10(1):110. Number: 1 Publisher: Nature Publishing	763
710	ciech Zaremba, Rowan Zellers, Chong Zhang, Mar-	Group.	764

Jingfeng Yang, Aditya Gupta, Shyam Upadhyay, Luheng He, Rahul Goel, and Shachi Paul. 2022. [TableFormer: Robust Transformer Modeling for Table-Text Encoding](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 528–537, Dublin, Ireland. Association for Computational Linguistics.

Xinyi Zheng, Doug Burdick, Lucian Popa, Xu Zhong, and Nancy Xin Ru Wang. 2020. [Global Table Extractor \(GTE\): A Framework for Joint Table Identification and Cell Structure Recognition Using Visual Context](#). ArXiv:2005.00589 [cs].

Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. 2020. [Image-based table recognition: data, model, and evaluation](#). ArXiv:1911.10683 [cs].

Yitan Zhu, Thomas Brettin, Fangfang Xia, Alexander Partin, Maulik Shukla, Hyunseung Yoo, Yvonne A. Evrard, James H. Doroshov, and Rick L. Stevens. 2021. [Converting tabular data into images for deep learning with convolutional neural networks](#). *Scientific Reports*, 11(1):11325. Number: 1 Publisher: Nature Publishing Group.

A ESG Keywords

Environment

1. Scope 1 GHG Emissions

Scope 1 are all direct emissions from the activities of an organization under their control. This includes fuel combustion on site such as gas boilers, fleet vehicles and air-conditioning leaks.

2. Scope 2 GHG Emissions Market Volume

Scope 2 are indirect emissions from electricity purchased and used by the organization. Emissions are created during the production of the energy and eventually used by the organization. A market-based method reflects emissions from electricity that companies have actively chosen to purchase or reflects their lack of choice.

3. Scope 2 GHG Emissions Location Volume

Scope 2 emissions are indirect emissions from the generation of purchased energy. A location-based method reflects the average emissions intensity of grids on which energy consumption occurs (using mostly grid-average emission factor data)

4. Scope 2 GHG Emissions Other Volume

Scope 2 emissions are indirect emissions from the generation of purchased energy. Overall, if not clearly defined whether it is market-based calculation or location-based calculation

5. Scope 3 GHG Emissions

Scope 3 emissions are all other indirect emissions (excluding Scope 2) that occur in the value chain of the reporting company, including both upstream and downstream emissions.

6. Environmental Restoration and Investment Initiatives Monetary Value

The fields represent the monetary value spent on environmental initiatives.

7. Total Water Discharged

The fields represent the overall volume of water discharged by a company.

8. Total Water Withdrawal

The fields represent the total volume of water withdrawn by a company.

9. Total Water Recycled

The fields represent the total volume of water recycled or reused by a company.

10. Toxic Air Emissions - NOx

The fields represent the total amount of nitrous oxide (NOx)emissions emitted by a company.

11. Toxic Air Emissions - SOx

The fields represent the total amount of sulfur oxide (Sox) emissions emitted by a company.

12. Toxic Air Emissions - Overall

The fields represent the total amount of air emissions emitted by a company.

13. Toxic Air Emissions - VOC

The fields represent the total amount of volatile organic compound (VOC) emissions emitted by the company.

14. Hazardous Waste - Disposed to Aquatic

The fields represent the total amount of hazardous waste disposed to aquatic environment.

15. Hazardous Waste - Disposed to Land

The fields represent the total amount of hazardous waste disposed to non aquatic or land environment.

16. Hazardous Waste - Total Recycled

The fields represent the total amount of hazardous waste recycled.

17. Hazardous Waste - Total Amount Generated

The fields represent the total amount of hazardous waste generated by a company.

18. Hazardous Waste - Total Amount Disposed

The fields represent the total amount of hazardous waste disposed.

19. Non-Hazardous Waste - Disposed to Aquatic

The fields represent the total amount of non-hazardous waste disposed to the aquatic environment.

870	20. Non-Hazardous Waste - Disposed to Land	33. Impacted Number of Species on National	922
871	The fields represent the total amount of non-	listed Species	923
872	hazardous waste to non aquatic or land envi-	The field identifies the number of impacted	924
873	ronment	species on National Listed Species.	925
874	21. Non-Hazardous Waste - Total Recycled	34. Baseline Level	926
875	The field represents the total amount of non-	The field identifies the value at baseline or	927
876	hazardous waste recycled.	year that target is set against.	928
877	22. Non-Hazardous Waste - Total Amount Gen-	35. Target Year	929
878	erated	The field identifies the year in which the re-	930
879	The fields represent the total amount of non-	newable energy goal is set to be completed.	931
880	hazardous waste Generated by a company.	36. Target Goal	932
881	23. Non-Hazardous Waste - Total Amount Dis-	The field identifies the target goal for renew-	933
882	posed	able energy.	934
883	The fields represent the total amount of non-	37. Actual Achieved	935
884	hazardous waste disposed.	The fields identifies the actual value achieved	936
885	24. Total Waste Produced	for the renewable energy goal.	937
886	The fields represent the total amount of waste	38. Baseline Level	938
887	produced by a company.	The field identifies the baseline emissions	939
888	25. Total Waste Recycled	value.	940
889	The fields represents the total amount of waste	39. Target Year	941
890	recycled by a company.	The field identifies the year in which GHG	942
891	26. Total Waste Disposed	emission goal is set to be completed.	943
892	This fields represent the total amount of waste	40. Target Goal	944
893	disposed by a company.	The field identifies the target goal for GHG	945
894	27. Number of Sites in Water Stress Areas	emission reduction.	946
895	The field represents the number of sites lo-	41. Actual Achieved	947
896	cated in water stress areas.	The field identifies the value achieved of GHG	948
897	28. E-Waste Produced	emissions reduced compare - in metric tons.	949
898	The field identifies the mass volume of f E-		
899	waste produced which are electronic products	Social	950
900	that are unwanted, not working, and nearing or	1. Training Hours Per Employee	951
901	at the end of their life. Examples of electronic	The fields identifies the numerical value of	952
902	waste include, but not limited to : computers,	training hours per employee.	953
903	printers, monitors, and mobile phones	2. Training Hours Annually	954
904	29. E-Waste Recycled	The fields identifies the numerical values of	955
905	The field identifies the mass volume of E-	training hours conducted within a year.	956
906	Waste Recycled.	3. Lost Time Injury Overall Rate	957
907	30. E-Waste Disposed	The fields identifies the total number of in-	958
908	The field identifies the mass volume of E-	juries that caused the employees and contrac-	959
909	waste disposed.	tors to lose at least a working day.	960
910	31. Number of Sites Operating in Protected	4. Lost Time Injury Rate Contractors	961
911	and/or High Biodiversity Areas	The fields identifies the number of injuries	962
912	The field identifies the number of sites or facil-	that caused the contractors to lose at least a	963
913	ities owned,leased, managed in or adjacent to	working day.	964
914	protected areas and areas of high biodiversity	5. Lost Time Injury Rate Employees	965
915	value outside protected areas.	The fields identifies the number of injuries	966
916	32. Impacted Number of Species on Interna-	that caused the employees to lose at least a	967
917	tional Union of Conservation of Nature	working day.	968
918	(IUCN) List	6. Employee Fatalities	969
919	The field identifies the number of impacted	The fields identifies the number of employee	970
920	species on International Union of Conserva-	fatalities during a one year period.	971
921	tion of Nature (IUCN) red list.	7. Contractor Fatalities	972
		The fields identifies the number of contractor	973
		fatalities during a one year period.	974

975	8. Public Fatalities		
976	The fields identifies the number of general		
977	public fatalities during a one year period.		
978	9. Number of Other Fatalities		
979	The fields identifies the number of fatalities		
980	during a one year period not broken down by		
981	employee, contractor, or public.		
982	10. Total Incident Rate Overall Workers		
983	The field identifies the number of work-related		
984	injuries per 100 overall workers during a one		
985	year period for both employees and contrac-		
986	tors.		
987	11. Total Incident Rate Contractors		
988	The field identifies the number of contractor		
989	work-related injuries per 100 overall workers		
990	during a one year period.		
991	12. Total Incident Rate Employees		
992	The field identifies the number of work-related		
993	injuries per 100 overall workers during a one		
994	year period for employees.		
995	13. Employee Turnover - Gender Male Rate		
996	The field identifies the absolute number		
997	turnover rate by males in a company .		
998	14. Employee Turnover - Gender Female Rate		
999	The field identifies the absolute number		
1000	turnover rate by females in a company.		
1001	15. Employee Turnover Overall Rate		
1002	The field identifies the absolute number		
1003	turnover rate for overall employees in a com-		
1004	pany.		
1005	16. Median Gender Pay Gap - Global		
1006	The field identifies the gender pay gap median		
1007	value of the company at a global level.		
1008	17. Mean Gender Pay Gap - Global		
1009	The field identifies the gender pay gap mean		
1010	or average value of the company at a global		
1011	level.		
1012	18. Median Gender Pay Gap by Location		
1013	The field represents the gender pay gap me-		
1014	dian value of the company at a location or		
1015	country level.		
1016	19. Mean Gender Pay Gap by Location		
1017	The field represents the gender pay gap		
1018	mean/average value of the company at a loca-		
1019	tion or country level.		
1020	20. Employee Turnover by Age - Lower Value		
1021	The field Identifies the minimum age in a		
1022	given range for employee turnover statistics.		
1023	21. Employee Turnover by Age - Upper Value		
1024	The field identifies the maximum age in a		
1025	given range for employee turnover statistics.		
1026	22. Employee Turnover by Age - Rate		
1027	The field identifies the employee turnover rate.		
	23. Employee Turnover by Location Rate	1028	
	The field identifies the absolute number of	1029	
	employee turnover rate by location.	1030	
	24. Workforce Breakdown Rate	1031	
	The field identifies the absolute number of	1032	
	employees of a company based on seniority,	1033	
	ethnicity or gender.	1034	
	25. Workforce Breakdown Job Category Data:	1035	
	Value (ABS)	1036	
	The field represents the employee count abso-	1037	
	lute value at a category level within a work-	1038	
	force.	1039	
	26. Number Of Product Recalls	1040	
	The fields identifies the number of product	1041	
	recalls.	1042	
	27. Product Recalls Annual Recall Rate	1043	
	The fields identifies the product recall rate of	1044	
	a company.	1045	
	Governance	1046	
	1. Percentage of Negative Votes on Pay	1047	
	Practices Year	1048	
		1049	
	2. Board of Director Term Limit	1050	
	The field identifies maximum amount of years	1051	
	a board member can serve.	1052	
	3. Board of Director Term Duration	1053	
	The field identifies number of years a board	1054	
	member can serve before reelection.	1055	
	4. Auditor Election Year	1056	
	The field identifies when the current lead au-	1057	
	ditor elected.	1058	
	5. Independent Auditor Start Year	1059	
	The field represents the start year the com-	1060	
	pany started having the audit company as its	1061	
	independent auditor.	1062	
	6. Average/Mean Compensation of Company	1063	
	Employees-Global	1064	
	The field represents the average or mean com-	1065	
	penetration for company employeesat a global	1066	
	level.	1067	
	7. Ratio Average Compensation of CEO to	1068	
	Employee - CEO- Global	1069	
	The field represents the ratio between the com-	1070	
	penetration paid to the companies CEO and the	1071	
	average compensations received by employ-	1072	
	ees at a global level.	1073	
	8. Compensation of Company Employees by	1074	
	Location	1075	
	The field identifies the average compensation	1076	
	for company employees at a location level.	1077	
	9. Number of Suppliers Complying with Code	1078	
	of Conduct	1079	

1080	The field identifies the number of suppliers	26. CEO Share Ownership	1133
1081	that comply with companies supplier code of	The field identifies the number of shares the	1134
1082	conduct.	CEO owns in the company.	1135
1083	10. Share Class Numeric	27. CEO Share Class Numeric	1136
1084	The field identifies the share class numeric	The field identifies the share class numeric	1137
1085	component.	component.	1138
1086	11. Voting Rights	28. Board Member Age	1139
1087	The field identifies the number of voting rights	The field identifies the age of the members of	1140
1088	per each share of stock within each class.	the board.	1141
1089	12. Shares Outstanding	29. Board Member Term in Years	1142
1090	The field identifies the number of shares out-	The fields identifies how long the individual	1143
1091	standing within a companies common stock.	board member has been on the board which is	1144
1092	13. Chairman Effective Begin Year	determined in years.	1145
1093	The field indicates the year when the current	30. Board Member Effective Year (Director	1146
1094	chairman assume his or her position. This	Since)	1147
1095	field is used if a full effective date is not avail-	The fields identifies the year the individual	1148
1096	able.	board member started serving on the board.	1149
1097	14. Chairman Effective End Year	31. Board Profile As of Year	1150
1098	The field indicates the year when the chairman	The field identifies the year of the board infor-	1151
1099	left the position.	mation. An example would be the year of the	1152
1100	15. CEO Effective Begin Year	proxy statement.	1153
1101	The field identifies the year the CEO assumed	32. Participation On Other Company Board	1154
1102	his or her position.	The field identifies the number of boards a	1155
1103	16. CEO Effective End Year	member is part of outside of the organization.	1156
1104	The field indicates the year when the CEO left	33. For Value Negative Votes on Directors	1157
1105	the position.	The field identifies the number of for value	1158
1106	17. CEO Compensation Salary	votes the director received.	1159
1107	The field identifies the current CEO salary.	34. Against Value Negative Votes on Directors	1160
1108	18. CEO Compensation Overall	The field identifies the number of against votes	1161
1109	The field identifies the CEO's overall com-	the director received.	1162
1110	compensation including salary, bonuses and all	35. Abstain Value Negative Votes on Directors	1163
1111	awards.	The field identifies the number of votes that	1164
1112	19. CEO Cash Bonus	were abstained for a given director.	1165
1113	The field identifies the cash bonus value for	36. Broker Non Vote Value Negative Votes on	1166
1114	the CEO.	Directors	1167
1115	20. CEO Stock Award Bonus	The field identifies the number of broker non	1168
1116	The CEO Stock Award Bonus value	votes for given director.	1169
1117	21. CEO Option Awards	37. Number of Board Meetings Attended by	1170
1118	The CEO Option Awards bonus value	Board Member	1171
1119	22. CEO Other Awards	The field identifies the number of board meet-	1172
1120	The fields identifies other compensation out-	ings attended by a board member.	1173
1121	side of salary, cash bonus, stock award bonus	38. Number of Board Meetings Held by Com-	1174
1122	and option awards. This could include change	pany	1175
1123	in pension and values categorized as "all other	The field identifies the number of board meet-	1176
1124	compensation"	ings held by a company while member was on	1177
1125	23. CEO Pension	the board.	1178
1126	The fields identifies the CEO pension amount.	39. Total Members on Board per Skill Set	1179
1127	24. Cash Severance Value	The field identifies the number of board mem-	1180
1128	The fields identifies the amount of cash the	bers within a specific skillset type.	1181
1129	severance policy for each category.		
1130	25. Total Severance Value		
1131	The fields identifies the total value amount of		
1132	the severance policy.		

B Algorithm for Statement Extraction

1182

We present the algorithm we used to extract statements. For this algorithm, the inputs are the original table and the labels table.

1183

1184

Algorithm 1 Extract Statements

```
1: procedure EXTRACT STATEMENTS(Table, LabelsTable)
2:   Input: Table, LabelsTable: Table and Table of cell annotations
3:   AllStatements  $\leftarrow$  empty list
4:   for all row in LabelsTable do
5:     for all column in LabelsTable do
6:       if LabelsTable[row][column] = Property Value then
7:         Search in the same row and column for (Sub)-Property
8:         if Property is found then
9:           Append Headers in hierarchy to Property, if any, starting from the minimum level
10:          Construct Statement with Property, Row and Column
11:        else if SubProperty is found then
12:          Append Property to the SubProperty
13:          Append Headers in hierarchy to SubProperty, if any, starting from the maximum level
14:          Construct Statement with SubProperty, Row and Column
15:        else
16:          Property is not found, continue to the next iteration
17:        end if
18:        Append Statement to AllStatements
19:      end if
20:    end for
21:  end for
22:  Return AllStatements
23: end procedure
```

```

1: procedure CONSTRUCT STATEMENT(Row, Column, Property)
2:   Input: Row, Column, Property: Row and Column of the Property Value, with its related Property
3:   Output: Statement: list
4:   Statement  $\leftarrow$  empty list
5:   Predicate  $\leftarrow$  empty dictionary
6:   Predicate [Property Value]  $\leftarrow$  Table[Row][Column]
7:   Predicate [Property]  $\leftarrow$  Property
8:   Search in the same row and column(Unit Value)
9:   Predicate[Unit]  $\leftarrow$  Table[rowuv][columnuv]
10:  Search for a Subject - Subject Value pair
11:  Predicate[Subject]  $\leftarrow$  Table[rows][columns]
12:  Predicate[Subject_Value]  $\leftarrow$  Table[rowsv][columnsv]
13:  Add Predicate to the Statement
14:  Search in the same row and column(Time Value)
15:  if Time Value is found then
16:    Predicate  $\leftarrow$  empty dictionary
17:    Predicate [Property Value]  $\leftarrow$  Table[rowtv][columntv]
18:    Predicate [Property]  $\leftarrow$  "Time"
19:    Add Predicate to the Statement
20:  end if
21:  Search for all Key - Key Value pairs
22:  for all Key - Key Value pairs found do
23:    Predicate  $\leftarrow$  empty dictionary
24:    Predicate[Property]  $\leftarrow$  Table[rowk][columnk]
25:    Predicate[Property Value]  $\leftarrow$  Table[rowkv][columnkv]
26:    Add Predicate to the Statement
27:  end for
28:  Return Statement
29: end procedure

```

```

1: procedure APPEND HEADERS(Row, Column, Property, Level)
2:   Input: Row, Column, Property, Level: Row, Column, value of a Property cell and the level of the header to search for.
3:   Output: Property: string
4:   for all Rowa above Row do
5:     for all Columnl on the left of Column do
6:       if LabelsTable[Rowa][Columnl] is a header with a higher level than Level then
7:         Append Table[Rowa][Columnl] on top of Property
8:         if the level of LabelsTable[Rowa][Columnl] is maximum then
9:           Return Property
10:        else
11:          Append Headers in hierarchy to Property starting from the level of LabelsTable[Rowa][Columnl]
12:          Return Property
13:        end if
14:      end if
15:    end for
16:  end for
17:  Return Property
18: end procedure

```

Algorithm 2 Utility Functions

```

1: procedure APPEND PROPERTY(Row, Column, SubProperty)
2:   Input: Row, Column, SubProperty: Row, Column and Value of a SubProperty cell
3:   Output: Subproperty: string
4:   for all Rowa above Row do
5:     for all Columnl on the left of Column do
6:       if LabelsTable[Rowa][Columnl] is a Property then
7:         Append Table[Rowa][Columnl] on top of SubProperty
8:         Return SubProperty
9:       end if
10:    end for
11:  end for
12:  Return SubProperty
13: end procedure

```

```

1: procedure SEARCH IN THE SAME ROW AND COLUMN(Row, Column, Key)
2:   Input: Row, Column, Key: Row and Column where to search the specified Key
3:   Output: Rowk, Columnk: Row and column of the designated Key, if found
4:   for all Cell respectively on the Left, Above, and Right to the cell at LabelsTable[Row][Column] do
5:     if Cell is Key then
6:       Return Row, Column of Cell
7:     end if
8:   end for
9:   Return Null
10: end procedure

```

Algorithm 3 Utility Functions

```

1: procedure SEARCH FOR A PAIR(Row, Column, Key, Key Value)
2:   Input: Row, Column, Key: Row and Column where to search the specified Key
3:   Output: Rowk, Columnk: Row and column of the designated Key, if found
4:   for all Cellkv respectively on the Left, Above, and Right to the cell at LabelsTable[Row][Column] do
5:     if Cellkv is Key Value then
6:       for all Cellk in the Orthogonal Direction with respect to Cellkv from LabelsTable[Row][Column] do
7:         if Cellk is Key then
8:           Return Coordinates of Cellk, Cellkv
9:         end if
10:      end for
11:    end if
12:  end for
13:  Return Null
14: end procedure

```
