# Entity Cloze By Date: What LMs Know About Unseen Entities

**Yasumasa Onoe, Michael J.Q. Zhang, Eunsol Choi, Greg Durrett**
Department of Computer Science
The University of Texas at Austin
`{yasumasa, mjqzhang, eunsol, gdurrett}@cs.utexas.edu`

## Abstract

Language models (LMs) are typically trained once on a large-scale corpus and used for years without being updated. However, in a dynamic world, new entities constantly arise. We propose a framework to analyze what LMs can infer about new entities that did not exist when the LMs were pretrained. We derive a dataset of entities indexed by their origination date and paired with their English Wikipedia articles, from which we can find sentences about each entity. We evaluate LMs' perplexity on masked spans within these sentences. We show that models more informed about the entities, such as those with access to a textual definition of them, achieve lower perplexity on this benchmark. Our experimental results demonstrate that making inferences about new entities remains difficult for LMs. Given its wide coverage on entity knowledge and temporal indexing, our dataset can be used to evaluate LMs and techniques designed to modify or extend their knowledge. Our automatic data collection pipeline can be easily used to continually update our benchmark.

## 1 Introduction

New entities arise every day: new movies, TV shows, and products are created, new events occur, and new people come into the spotlight. Whatever the capabilities of language models (LMs) to represent entity knowledge, these new entities cannot possibly be included in the language models' parametric knowledge (i.e., knowledge acquired during pretraining), as they did not exist when LMs were trained. Since this temporal mismatch between LMs and real-world knowledge affects model performance on downstream tasks (Zhang and Choi, 2021; Dhingra et al., 2021; Lazaridou et al., 2021), understanding what LMs know about real-world entities is an important task.

The existing literature provides various benchmarks to measure LMs' knowledge about entities (Petroni et al., 2019, 2021; Dhingra et al.,
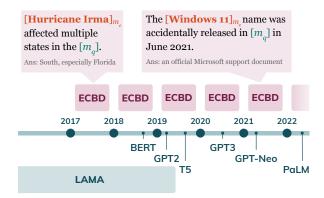


Figure 1: Our framework (ECBD) collects entities indexed by the year when they were first introduced in Wikipedia and their cloze sentences, unlike existing cloze datasets (LAMA (Petroni et al., 2019)) which broadly cover entities introduced prior to 2019.

2021). Those benchmarks are typically formulated as cloze-style tasks covering a limited set of relations bounded by knowledge bases: LAMA uses around 40 Wikidata relations and entities collected in 2017. Newer cloze benchmarks (Dhingra et al., 2021; Jang et al., 2021) integrate temporal aspects to identify a time period when a cloze sentence is valid, but do not differentiate new and existing entities. These knowledge probing datasets fail to test broad knowledge about real-world entities or evaluate how LMs' knowledge differs on entities that are seen or unseen during pre-training.

To fill this gap, we propose a framework to evaluate LMs' knowledge about entities classified by their origination date. We extract a set of Origination Date Indexed Entities (ODIE) based on metadata from Wikidata. We then construct cloze statements by masking sentences in those entities' Wikipedia articles. Unlike past knowledge probing datasets, these cloze sentences test the ability of a model to make a wide range of inferences related to entities, without being restricted to a pre-defined set of KB relations. We choose masked spans near these entities that likely contain information related to the entities, which we evaluate based on

the perplexity gap between the raw sentence and the sentence with the entity replaced.

We release the Entity Cloze by Date (ECBD) dataset of 35k masked sentences that contain mentions of 2.1K ODIE entities,[1] split by year covering a time period from 2017 to 2021, together with 8k masked sentences of popular entities from any time period. In our experiments, we evaluate three pre-trained language models in terms of perplexity. We establish that injecting additional information such as a text definition can meaningfully teach the model to make better guesses about masked spans, highlighting this dataset's utility for benchmarking methods of knowledge injection.

## 2    Entity Cloze by Date

We aim to test language models' 1) broader entity knowledge and 2) ability to reason about completely unseen entities (i.e., unseen during pretraining). Thus, we want to have the following properties in our entity cloze sentences. **(1) Date indexing.** If each cloze example is associated with an entity and indexed by the origination date of that entity, we can understand whether a model may have seen it in its pre-training corpus or not. **(2) Diverse sentences.** When going beyond KB triples, entity knowledge can take many forms: actions that an entity can take, other entities that action can effect, typical ways in which an entity is described, and more. Thus, we want include diverse sentences and masked spans that cover rich relations and various syntactic categories (e.g., POS and nonterminal categories, span length).

### 2.1    Task Definition

Each entity $e$ is paired with $e_i$, its origination year. Given a sentence $s$ containing an entity mention span $m_e$ and a masked query span $m_q$, a language model is asked to predict the gold masked span $m_y$. See the following example:

$e$: RNA vaccine, $e_i$: 2020
$s$: [mRNA vaccines]$_{m_e}$ do not affect or reprogram [$m_q$].
$m_y$: DNA inside the cell

We evaluate language models by **perplexity** on the masked span $m_q$ (see Appendix D for a discussion of recall as another metric).

Figure 2: Overview of the data collection process.

### 2.2    Data Collection

Our data collection protocol consists of three stages: *entity mining*, *sentence collection* and *span selection*. We use English Wikipedia (the September 1, 2021 dump) and Wikidata as knowledge sources.

**ODIE Mining**    We begin by gathering all entities on Wikidata that have an associated *start time*, *announcement date*, *time of discovery or invention*, *inception* date, *point in time*, or date it was *introduced on*. For such entities, we take the first of these dates to create our temporal splits, assuming that this is the earliest date the entity could have appeared in any pretraining corpus.

To compare with ODIE which covers relatively new entities originated in 2017 at the earliest, we use a set of POPULAR entities ranked by article contributor numbers and incoming links from prior work (Onoe et al., 2021; Geva et al., 2021).

**Entity Sentence Collection**    Once we obtain a list of entities, we look up their English Wikipedia articles. To enrich the candidate sentence pool and exclude trivial sentences from stub articles, we filter entities if their corresponding articles contain less than 500 words. From each article, we exclude the first paragraph of the article, to be used as an entity definition, and sample sentences from the rest of the paragraphs. We sample sentences that include the entity name or one of their Wikidata aliases. We do not accept entity mention spans located in quotes since they are often in nested named entities such as book titles. We also filter out any sentences with less than five words.

**Span Selection**    Next, we determine spans $m_q$ to be masked on a sentence, $s$; we can have multiple masked spans per sentence, masked separately. All spans must be: (a) not overlapping with the entity mention span, $m_e$, (b) located after the entity mention span, $m_e$, and (c) starting no more than ten words away from the mention span, to improve relatedness to the entity. We select spans after the entity mention so left-to-right language models will

| Origination Year | 2017 | 2018 | 2019 | 2020 | 2021 | Total | Example Entities |
|---|---|---|---|---|---|---|---|
| # Dev Entities | 300 | 280 | 219 | 187 | 78 | 1,050 | |
| # Test Entities | 299 | 279 | 208 | 176 | 80 | 1,029 | |
| Sports | 20 | 19 | 22 | 12 | 27 | 19 | 2017 Tour de France, USL League One, Evo 2017 |
| Media | 18 | 19 | 24 | 23 | 20 | 21 | Emily in Paris, Luigi's Mansion 3, The Midnight Gospel |
| Infrastructure | 10 | 8 | 10 | 8 | 9 | 9 | Gateway Arch National Park, Istanbul Airport, I-74 Bridge |
| Natural Risks | 3 | 6 | 4 | 15 | 11 | 7 | Hurricane Ida, COVID-19, North Complex Fire |
| Products | 4 | 4 | 4 | 3 | 3 | 4 | Apple Card, Sputnik V COVID-19 vaccine, Pixel 4 |
| Businesses | 15 | 11 | 7 | 7 | 3 | 10 | Raytheon Technologies, Electrify America, Good Party |
| Organizations | 16 | 18 | 13 | 12 | 9 | 15 | NUMTOT, UK Student Climate Network |
| Other Events | 9 | 10 | 11 | 12 | 13 | 11 | Super Bowl LIV halftime show, Storm Area 51 |
| Misc. | 5 | 3 | 4 | 7 | 4 | 4 | RNA vaccine, Earthshot Prize, Comet NEOWISE |

Table 1: Origination date indexed entity (ODIE) statistics by category. The number represents % of entities with particular type among entities originated in that year.

condition on the entity at test time.

We extract two types of spans: **NP spans** are selected from any suitable noun phrases in the sentence using spaCy (Honnibal and Montani, 2017). These spans primarily represent relational knowledge about the entity, analogous to the object in a KB triple. **Random spans** are arbitrary sequences of words sampled from the sentence. This broader set of spans may cover other types of entity knowledge (e.g., probable actions an entity can take). We uniformly sample span length between 1 and 5 and then randomly select the starting location of the span within the sentence. We only accept valid spans not overlapping with the entity mention. We extract at most 100 spans per entity to limit any one entity's contribution to the final dataset.

**Span sensitivity to entity knowledge**   To see if our design choices are effective, we perform a test that measures the performance drop in perplexity using T5 when we replace the entity mention with a generic reference to "the entity." We use entities from our POPULAR set to ensure that the LM has seen them during pre-training. If a masked span is related to the entity, the perplexity of that span should increase when the entity mention is omitted. We see that the median perplexity of a span increases by 32.2% when the entity is removed, indicating that these spans are indeed related to the entity. Moreover, removing the distance-based criterion for span selection decreases the perplexity change to 25.9%. These results indicate that our selected spans are correlated with the entity. This gap test was performed only for analysis and we do not use any model-based data filtering.

**Dataset Statistics**   Table 1 shows the statistics and examples of ODIE, split by entity types. While our entity set does not comprehensively capture all entities originating in that year, it contains a diverse

| | # Sent. | # Ent. | $m_q$ Span Len. | |Span V.| |
|---|---|---|---|---|
| LAMA$_{TREx}$ | 34k | 29,488 | 1.0 | 2,017 |
| ECBD | 35k | 2,106 | 2.9 | 19,542 |
| POPULAR | 8k | 1,910 | 2.9 | 8,644 |

Table 2: Data statistics. |Span V.| means the vocabulary size of masked spans. Initial release of the data sample equal number of masked sentences per year (2017-21).

set of entities, ranging from events, products to organizations. One notably missing entity category is people; it is hard to pin down an origination year because of the significant gap between birth year and the year someone became prominent.

Table 2 reports statistics on our cloze task data and existing probe dataset (Petroni et al., 2019). While containing fewer entities, our dataset exhibits much richer vocabulary (19K vs. 2K), demonstrating diverse knowledge it covers. We split this data into dev and test sets by entities (i.e., no shared entities between dev and test). To balance out the data sizes across the groups, we sample 4k examples for each year group, yielding 35k examples in total (approx. 20k for dev and 20k for test). Earlier dates contain a larger set of entities (599 entities for 2017 compared to 158 entities for 2021) as entities are continuously updated in Wikidata. In other words, many entities originated in 2021 have not been yet added to Wikidata. We sample the same number of NP spans and random spans. Within the NP spans, 35% of them are proper noun phrases.

## 3   Experiments

**Setup**   We evaluate T5-large (Raffel et al., 2020), BART-large (Lewis et al., 2020), and GPT-Neo (Black et al., 2021) on our dataset in the zero-shot setting where the model parameters are fixed. In addition to the original masked sentence

COVAX began distributing `<extra_id_0>` in February 2021. ➡️ Compute token-normalized perplexity

`<extra_id_0>` _vaccine s `<extra_id_1>`...

*Input Sequence: Left and Right Context*     *Target Sequence*

**BART** (seq-to-seq, BPE tokenizer)

COVAX began distributing `<mask>` in February 2021. ➡️ Compute token-normalized perplexity

COVAX began distributing vaccines in February 2021.

*Input Sequence: Left and Right Context*     *Target Sequence*

**GPT-Neo** (left-to-right, BPE tokenizer)

GPT-Neo only receives the left context as input.      Compute token-normalized perplexity

COVAX began distributing vaccines in February 2021.
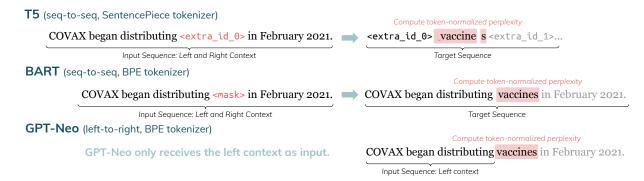
*Input Sequence: Left context*

Figure 3: Perplexity computation over the masked span with three different modeling paradigms.

(ORIGINAL), we feed three modified masked sentences. NO ENT replaces the entity mention span with a generic string "the entity." RANDOM DEF. prepends a definition sentence of a randomly selected entity. DEFINITION prepends the first sentence of the entity's Wikipedia article to the cloze sentence.

We evaluate these models on the subsets split by year as well as a set of popular entities. Note that the entities in the 2020 and 2021 subsets are unseen for T5 and BART. Most entities in the 2020 and 2021 subsets are unseen to GPT-Neo, but its training data (the Pile (Gao et al., 2020)) does include the March 2020 English Wikipedia dump. In our experiments, we group the 2020 and 2021 subsets together as they consist of "unseen" entities. Similarly, we group the 2017, 2018, and 2019 subsets whose entities are "seen" during pre-training. See Appendix B for perplexity per year.

**Evaluation Metric** We compute token-normalized perplexity over the span as a proxy for entity knowledge stored in LMs. Each subset has different distribution of entity types (e.g., 2020 contains many COVID related entities and a lot less sports events compared to other years), and some frequent entities might contribute to perplexity excessively. To mitigate biases from particular entities, we first average negative log-likelihood (token normalized) over entities then average over examples. We follow the target sequence format used in LMs' pre-training tasks (see Figure 3).

Figure 3 shows the perplexity computation. For left-to-right language models like GPT-Neo, we compute the perplexity of the span given the left context *only*. T5 and BART, as seq2seq models, are able to also condition on the right context in their input; this makes perplexity values between these model classes not directly comparable (in addi-

| | POPULAR | 2017-2019 | 2020-2021 |
|---|---|---|---|
| Type: seq-to-seq | **T5 Large** | | Size: 770M |
| ORIGINAL | 13.02 | 15.39 | 19.43 |
| NO ENT | 18.28 | 22.35 | 26.69 |
| RANDOM DEF. | 12.10 | 14.33 | 17.34 |
| DEFINITION | 11.04 | 11.73 | 13.60 |
| Δ (ORIG. → RAND.) | -0.92 | -1.06 | -2.09 |
| Δ (ORIG. → DEF.) | -1.98 | -3.66 | -5.83 |
| Type: seq-to-seq | **BART Large** | | Size: 406M |
| ORIGINAL | 22.70 | 21.09 | 28.79 |
| NO ENT | 33.33 | 30.56 | 39.25 |
| RANDOM DEF. | 27.69 | 25.59 | 33.74 |
| DEFINITION | 21.10 | 17.66 | 22.00 |
| Δ (ORIG. → RAND.) | +4.99 | +4.50 | +4.95 |
| Δ (ORIG. → DEF.) | -1.60 | -3.43 | -6.79 |
| Type: left-to-right | **GPT-Neo** | | Size: 1.3B |
| ORIGINAL | 28.61 | 27.81 | 33.36 |
| NO ENT | 54.01 | 51.46 | 54.81 |
| RANDOM DEF. | 39.46 | 41.03 | 45.92 |
| DEFINITION | 23.19 | 19.09 | 22.33 |
| Δ (ORIG. → RAND.) | +10.85 | +13.22 | +12.56 |
| Δ (ORIG. → DEF.) | -5.42 | -8.72 | -11.03 |

Table 3: Results of T5, BART, and GPT-Neo on the test set, showing perplexity (↓).

tion to differences in tokenization and pre-training tasks). For T5 and BART, we condition on the input with a single mask. At decoding, for BART we initialize the decoder with the left context of the span and compute perplexity on the true span filler following this left context. For T5, we compute perplexity on the output span between the special tokens `<extra_id_0>` and `<extra_id_1>`.

**Results** Table 3 reports perplexity (lower is better) on the test set that is split into three subsets: POPULAR, 2017-2019, and 2020-2021. Note that absolute perplexity across years is sensitive to factors such as distribution of sentences or entity types; we thus focus on relative performance.

In all subsets, we observe two consistent trends across three LMs. (1) NO ENT always degrades

performance compared to ORIGINAL. This result confirms that our masked spans are sensitive to the content of the entity span, although it is not conclusive proof of entity knowledge being required, as changing to "the entity" modifies other latent stylistic attributes that the LMs may be sensitive to. (2) DEFINITION always boosts performance over ORIGINAL, indicating that providing more information about entities helps to retrieve information distributed over LMs' parameters. RANDOM DEF. distracts BART and GPT-Neo but slightly improves T5 performance even though the additional information is taken from a random entity. This could be due to the model using different positional encodings as a result of having a definition, or LMs may select information if it is useful in some cases, leading the small gains.

**Performance on unseen entities** Recall that we consider 2020-2021 as unseen entities, and 2017-2019 and POPULAR as seen entities. All three LMs give higher perplexity on unseen entities, showing that the spans in 2020-2021 are relatively unexpected to the LMs.

We further investigate the performance delta between ORIGINAL and DEFINITION per subset. For all three LMs, we see that the performace delta is relatively larger on 2020-2021, indicating definition sentences are more useful on unseen entities.

Also, the performance delta on the popular entity set is notably smaller than 2020-2021 (compare T5 numbers: $13.02 \rightarrow 11.04$ for POPULAR versus $19.43 \rightarrow 13.60$ for 2020-2021). This implies that LMs contain some prior knowledge about common entities they have observed before, and can use additional information about new entities or less frequent entities. How to inject knowledge requires further investigation.

## 4 Use Cases

We envision this dataset as being useful for general knowledge probing, as the real-world knowledge covered by the existing benchmarks is gradually outdated. With our framework, we can easily **update** datasets using the most recent knowledge sources with a controlled manner. Since the entity knowledge in our dataset is time-indexed, this is suitable for evaluating knowledge editing approaches (Sinitsin et al., 2020; Zhu et al., 2020; De Cao et al., 2021; Mitchell et al., 2021; Meng et al., 2022) and also continual knowledge learning approaches (Jang et al., 2021). Crucially, existing work studies whether these approaches can inject single facts, but not whether they can enable models to robustly support a broad range of new inferences about entities, like our dataset allows.

## 5 Related Work

Temporal mismatch/misalignment between large pre-trained LMs and real-world knowledge is an emerging research direction. Lazaridou et al. (2021) show that the corpus-level perplexity on documents from beyond LMs' training period becomes increasingly poor over time. Dhingra et al. (2021) propose TEMPORALLAMA, which is based on time-dependent knowledge base triples (i.e., valid subject, relation, and object combinations given time). SITUATEDQA (Zhang and Choi, 2021) includes time-dependent QA examples. While these datasets primarily test temporal information about entities in the pre-training data, ECBD focuses on new entities which did not exist during pre-training. TemporalWiki (Jang et al., 2022) annotates new facts/entities based on the differences between Wikidata/English Wikipedia dumps, but does not necessarily reflect real-world changes during the time period (e.g., an ancient queen can be added to Wikidata in 2022). ECBD selects entities based on their origination date to align them with the real-world timeline.

Another line of work has looked at diachronic embeddings: (Wijaya and Yeniterzi, 2011; Kim et al., 2014; Hamilton et al., 2016; Bamler and Mandt, 2017), which can model changing meanings of words over time. Our setting focuses on introducing *new* concepts rather than rewriting existing ones, but data similar to ECBD could be collected for new usages of existing words.

Although our dataset follows the widely-used cloze format, our focus is orthogonal to datasets like the Children's Book Test (Hill et al., 2016) and LAMBADA (Paperno et al., 2016), which come from fiction and do not cover real-world entities.

## 6 Conclusion

In this paper, we present a dataset to understand language models' broad inferences about entities across time. We collect 43k cloze-style sentences associated with a time-indexed set of entities. We also perform analysis on our data set and show that handling completely unseen entities remains challenging for the current LMs.

## Acknowledgments

## References

Robert Bamler and Stephan Mandt. 2017. Dynamic Word Embeddings. In *Proceedings of the International Conference on Machine Learning (ICML)*.

Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. If you use this software, please cite it using these metadata.

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Bhuwan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W. Cohen. 2021. Time-Aware Language Models as Temporal Knowledge Bases.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361.

William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.

Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2016. The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations. In *International Conference on Learning Representations (ICLR)*.

Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.

Joel Jang, Seonghyeon Ye, Changho Lee, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, and Minjoon Seo. 2022. TemporalWiki: A Lifelong Benchmark for Training and Evaluating Ever-Evolving Language Models. abs/2204.14211.

Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Jungkyu Choi, and Minjoon Seo. 2021. Towards Continual Knowledge Learning of Language Models.

Yoon Kim, Yi-I Chiu, Kentaro Hanaki, Darshan Hegde, and Slav Petrov. 2014. Temporal Analysis of Language through Neural Language Models. In *Proceedings of the ACL 2014 Workshop on Language Technologies and Computational Social Science*.

Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d'Autume, Tomas Kocisky, Sebastian Ruder, Dani Yogatama, Kris Cao, Susannah Young, and Phil Blunsom. 2021. Mind the Gap: Assessing Temporal Generalization in Neural Language Models. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pretraining for natural language generation, translation, and comprehension. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual knowledge in gpt. *arXiv preprint arXiv:2202.05262*.

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2021. Fast Model Editing at Scale. abs/2110.11309.

Yasumasa Onoe, Michael J.Q. Zhang, Eunsol Choi, and Greg Durrett. 2021. CREAK: A Dataset for Commonsense Reasoning over Entity Knowledge. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track*.

Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and R. Fernández. 2016. The lambada dataset: Word prediction requiring a broad discourse context. *ArXiv*, abs/1606.06031.

Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean

Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. KILT: a benchmark for knowledge intensive language tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544, Online. Association for Computational Linguistics.

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)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.

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):1–67.

Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitriy Pyrkin, Sergei Popov, and Artem Babenko. 2020. Editable Neural Networks. In *International Conference on Learning Representations (ICLR)*.

Derry Tanti Wijaya and Reyyan Yeniterzi. 2011. Understanding semantic change of words over centuries. In *International Workshop on Detecting and Exploiting Cultural Diversity on the Social Web, DE-TECT*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Michael J.Q. Zhang and Eunsol Choi. 2021. SituatedQA: Incorporating extra-linguistic contexts into QA. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7371–7387, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv Kumar. 2020. Modifying memories in transformer models. *arXiv preprint arXiv:2012.00363*.

## A   Examples of ECBD Sentences

See Table 4 for examples of masked sentences in the ECBD data.

## B   Perplexity per year

See Table 5 for a more fine-grained view of the results in Table 3.

## C   Perplexity per span type

See Table 6 for a breakdown of the perplexity that T5 achieves on different types of spans, showing that random spans are generally higher perplexity than NP spans but that adding definitions can help both.

## D   Recall@10

LMs can be evaluated on **recall@10**, i.e., a binary score indicating if model's top ten predictions contains the gold masked span $m_y$. For T5, we first generate sequences using beam search (we choose beam size = 100 in our experiments). Then we take the top ten unique sequences and extract the text spans between `<extra_id_0>` and `<extra_id_1>` as predictions. Table 7 reports recall@10 on each subset. Table 8 list recall@10 per span type for each subset.

We only explore recall on T5, since it is not obvious how to compute it for the other two models. For BART, we can extract the predicted span by aligning the model's prediction with the gold context, assuming that it starts to copy from the input right context at some point. However, in some cases, we found that the generated right context does not match with the gold right context; it's unclear how to be handle this. For GPT-Neo, since it is a left-to-right LM, extracting the predicted span would require conditioning on the span length, which is information that T5 does not have access to. As a result, we do not report recall@10 for these models.

## E   Data Licensing

The Wikipedia text we used is licensed under CC BY-SA. Our use of Wikipedia, constructing a dataset which we will make publicly available under the same license, is consistent with the terms of the license.

## F   Computational Resources

All experiments were conducted using an NVIDIA Quadro RTX 8000. We only evaluate existing mod-

| Masked Sentence | Span Type | Origin Year |
|---|---|---|
| At 18:00 UTC on August 16, after **Grace** exited the Dominican Republic, [MASK] were lifted. Answer: "all tropical storm watches" | NP | 2021 |
| **AirTags** can be [MASK] the Find My app. Answer: "interacted with using" | RANDOM | 2021 |
| British tabloid "The Sun" is credited with the first headline use of '**Megxit**' on [MASK] 2020. Answer: "9 January" | NP | 2020 |
| The **iPhone SE** features an [MASK] a glass front and back. Answer: "aluminum frame, paired with" | RANDOM | 2020 |
| The **GPT-2** model has [MASK], and was trained on a dataset of 8 million web pages. Answer: "1.5 billion parameters" | NP | 2019 |
| The epicenter of the **2019 Albania earthquake** [MASK] kilometers from Tirana to the Northwest. Answer: "was about 30" | RANDOM | 2019 |
| On November 12, 2019, **Maverick City Music** released [MASK], "Maverick City, Vol. 2". Answer: "their follow-up EP" | NP | 2018 |
| **Austin FC** are the operators of a newly-[MASK]. Answer: "built stadium at McKalla Place" | RANDOM | 2018 |
| The first quarter of **Super Bowl LI** was [MASK] with each team punting twice. Answer: "a scoreless defensive match" | NP | 2017 |
| **Hurricane Irma** was the top Google searched term in [MASK] in 2017. Answer: "the U.S. and globally" | RANDOM | 2017 |

Table 4: Examples selected from the 2017-2021 subsets of ECBD.

els on our datasets and did not do any finetuning. One evaluation experiment typically takes 15 minutes to complete. For T5 experiments, we use Hugging Face's Transformer package (Wolf et al., 2020).

|  | POPULAR | 2017 | 2018 | 2019 | 2020 | 2021 |
|---|---|---|---|---|---|---|
| Type: seq-to-seq | | **T5 Large** | | | | Size: 770M |
| ORIGINAL | 13.02 | 15.28 | 14.78 | 16.43 | 19.81 | 18.60 |
| NO ENT | 18.28 | 22.28 | 21.70 | 23.35 | 28.41 | 23.26 |
| RANDOM DEF. | 12.10 | 14.56 | 13.54 | 15.10 | 17.42 | 17.17 |
| DEFINITION | 11.04 | 12.27 | 10.76 | 12.34 | 14.07 | 12.61 |
| Δ(ORIG. → RAND.) | -0.92 | -0.72 | -1.24 | -1.33 | -2.39 | -1.43 |
| Δ(ORIG. → DEF.) | -1.98 | -3.01 | -4.02 | -4.09 | -5.74 | -5.99 |
| Type: seq-to-seq | | **BART Large** | | | | Size: 406M |
| ORIGINAL | 22.70 | 22.74 | 19.52 | 21.00 | 28.03 | 30.53 |
| NO ENT | 33.33 | 33.58 | 28.25 | 29.67 | 39.56 | 38.57 |
| RANDOM DEF. | 27.69 | 27.11 | 23.80 | 25.96 | 32.41 | 36.86 |
| DEFINITION | 21.01 | 18.97 | 16.58 | 17.35 | 22.12 | 21.72 |
| Δ(ORIG. → RAND.) | +4.99 | +4.37 | +4.28 | +4.96 | +4.38 | +6.33 |
| Δ(ORIG. → DEF.) | -1.69 | -3.77 | -2.94 | -3.65 | -5.91 | -8.81 |
| Type: left-to-right | | **GPT-Neo** | | | | Size: 1.3B |
| ORIGINAL | 28.61 | 28.91 | 27.55 | 26.63 | 33.15 | 33.81 |
| NO ENT | 54.01 | 52.88 | 53.95 | 46.44 | 53.89 | 57.61 |
| RANDOM DEF. | 39.46 | 41.75 | 43.15 | 37.41 | 45.30 | 47.32 |
| DEFINITION | 23.19 | 20.47 | 18.00 | 18.68 | 22.17 | 22.69 |
| Δ(ORIG. → RAND.) | +10.85 | +12.84 | +15.6 | +10.78 | +12.15 | +13.51 |
| Δ(ORIG. → DEF.) | -5.42 | -8.44 | -9.55 | -7.95 | -10.98 | -11.12 |

Table 5: Results of T5, BART, and GPT-Neo on the test set, showing perplexity (↓) for each subset.

|  | 2017 | | 2018 | | 2019 | | 2020 | | 2021 | |
|---|---|---|---|---|---|---|---|---|---|---|
| Input Type | NP | RAND | NP | RAND | NP | RAND | NP | RAND | NP | RAND |
| ORIGINAL | 5.86 | 7.33 | 5.81 | 7.51 | 6.11 | 7.29 | 5.92 | 7.63 | 6.23 | 7.31 |
| NO ENT | 5.90 | 8.02 | 5.78 | 8.56 | 5.99 | 8.31 | 6.75 | 9.36 | 7.28 | 9.21 |
| RANDOM DEF. | 5.59 | 6.60 | 5.54 | 6.84 | 5.77 | 6.60 | 5.70 | 6.98 | 6.01 | 6.65 |
| DEFINITION | 4.96 | 5.98 | 4.98 | 6.02 | 5.12 | 5.85 | 5.14 | 6.13 | 5.13 | 5.82 |
| Δ(ORIG. → RAND.) | -0.27 | -0.73 | -0.27 | -0.67 | -0.34 | -0.69 | -0.22 | -0.65 | -0.22 | -0.66 |
| Δ(ORIG. → DEF.) | -0.90 | -1.35 | -0.83 | -1.49 | -0.99 | -1.44 | -0.78 | -1.50 | -1.10 | -1.49 |

Table 6: Results of T5 model (pre-trained with data from 2019) on the dev set with perplexity (↓) per span type.

|  | POPULAR | 2017 | 2018 | 2019 | 2020 | 2021 |
|---|---|---|---|---|---|---|
| Type: seq-to-seq | | **T5 Large** | | | | Size: 770M |
| ORIGINAL | 28.2 | 25.4 | 27.4 | 27.7 | 20.8 | 23.0 |
| NO ENT | 23.8 | 21.6 | 23.2 | 23.7 | 19.5 | 21.5 |
| RANDOM DEF. | 28.4 | 24.3 | 28.5 | 26.8 | 21.4 | 23.2 |
| DEFINITION | 29.3 | 28.4 | 31.8 | 28.2 | 24.8 | 26.1 |
| Δ(ORIG. → RAND.) | +0.2 | -1.1 | +1.1 | -0.9 | +0.6 | +0.2 |
| Δ(ORIG. → DEF.) | +1.1 | +3.0 | +4.4 | +0.5 | +4.0 | +3.1 |

Table 7: Results of T5 on the test set, showing recall@10 (↑) for each subset.

| Input Type | 2017 | | 2018 | | 2019 | | 2020 | | 2021 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | NP | Rand | NP | Rand | NP | Rand | NP | Rand | NP | Rand |
| Type: seq-to-seq | **T5 Large** | | | | | | | | | Size: 770M |
| Original | 30.3 | 20.0 | 31.8 | 20.2 | 29.3 | 22.0 | 30.1 | 19.8 | 29.3 | 21.6 |
| No Ent | 27.2 | 18.8 | 28.1 | 16.7 | 26.2 | 18.1 | 26.8 | 16.7 | 25.9 | 18.2 |
| Random Def. | 31.8 | 20.8 | 32.8 | 19.9 | 29.8 | 21.6 | 31.3 | 20.5 | 29.5 | 21.6 |
| Definition | 34.1 | 22.8 | 35.9 | 22.8 | 33.0 | 24.9 | 33.7 | 23.0 | 32.7 | 25.2 |
| $\Delta$(Orig. $\rightarrow$ Rand.) | +1.5 | +0.8 | +1.0 | -0.3 | +0.5 | -0.4 | +1.2 | +0.7 | +0.2 | +0.0 |
| $\Delta$(Orig. $\rightarrow$ Def.) | +3.8 | +2.8 | +4.1 | +2.6 | +3.7 | +2.9 | +3.6 | +3.2 | +3.4 | +3.6 |

Table 8: Results of T5 model (pre-trained with data from 2019) on the dev set with recall@10 (↑) per span type.