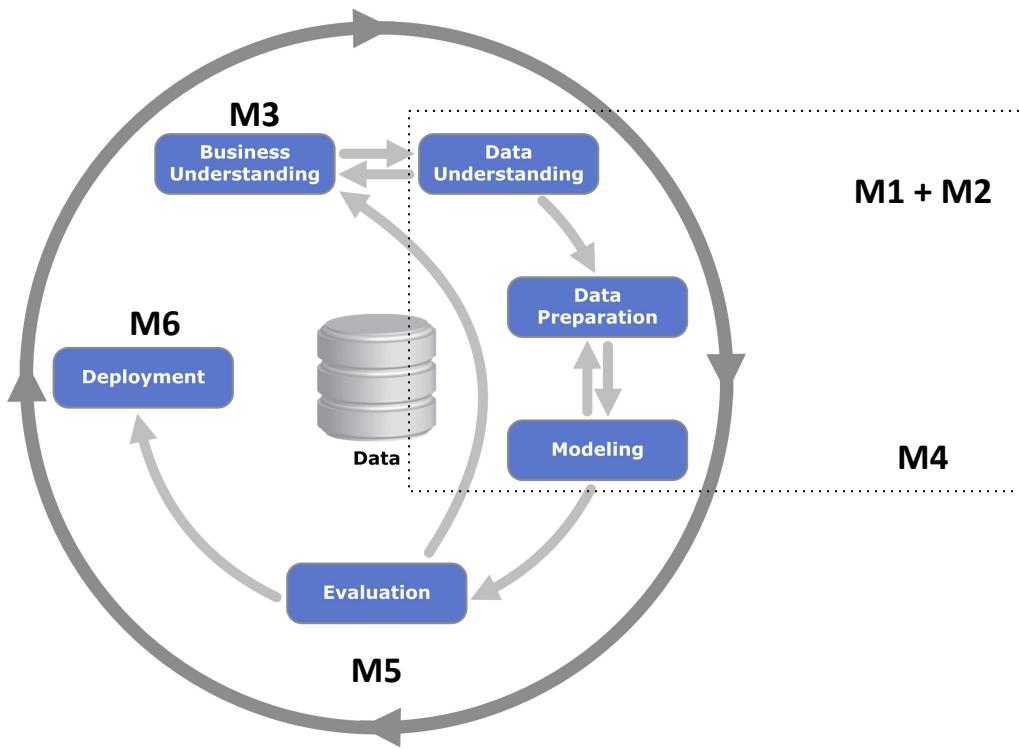


# Life2Vec, Market2Vec & Fine-Tuning

Hamid & Milad

Aalborg University Business School  
Innovation, Knowledge, and Economic Development (IKE)  
AI Denmark (AIDK)  
[hamidb@business.aau.dk](mailto:hamidb@business.aau.dk)  
Feb 02, 2026

# CRISP - DM



# Today's lecture overview

## Parts I & II

- Recap lecture 1
- Review the highlights in the BERT paper (Devlin et al. 2019)
- Fine tuning (API and Native Pytorch)

## Part III

- Socioeconomic phenomena as sequential processes
- Life2Vec
- Transformers and BDS
- Market2Vec
- Sentence Transformers Finetuning (SetFit)
- Going beyond Market2Vec

## Part IV

- Hands-on exploration: running code to understand how things work in practice

# Recap: Lecture 1

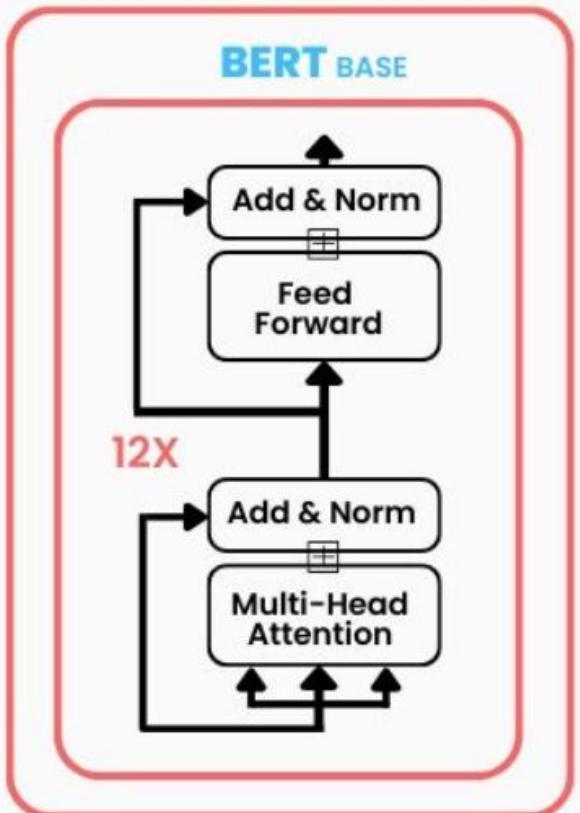
Let's refresh what we learned in the first lecture!

Attention and Contextualization

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova  
Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com



110M Parameters

arXiv:1810.04805v2 [cs.CL] 24 May 2019

## Abstract

We introduce a new language representation model called **BERT**, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications.

BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5% (7.7% point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute improvement) and SQuAD v2.0 Test F1 to 83.1 (5.1 point absolute improvement).

## 1 Introduction

Language model pre-training has been shown to be effective for improving many natural language processing tasks (Dai and Le, 2015; Peters et al., 2018a; Radford et al., 2018; Howard and Ruder, 2018). These include sentence-level tasks such as natural language inference (Bowman et al., 2015; Williams et al., 2018) and paraphrasing (Dolan

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as additional features. The fine-tuning approach, such as the Generative Pre-trained Transformer (OpenAI GPT) (Radford et al., 2018), introduces minimal task-specific parameters, and is trained on the downstream tasks by simply fine-tuning *all* pre-trained parameters. The two approaches share the same objective function during pre-training, where they use unidirectional language models to learn general language representations.

We argue that current techniques restrict the power of the pre-trained representations, especially for the fine-tuning approaches. The major limitation is that standard language models are unidirectional, and this limits the choice of architectures that can be used during pre-training. For example, in OpenAI GPT, the authors use a left-to-right architecture, where every token can only attend to previous tokens in the self-attention layers of the Transformer (Vaswani et al., 2017). Such restrictions are sub-optimal for sentence-level tasks, and could be very harmful when applying fine-tuning based approaches to token-level tasks such as question answering, where it is crucial to incorporate context from both directions.

In this paper, we improve the fine-tuning based approaches by proposing **BERT**: Bidirectional Encoder Representations from Transformers.

Comparison	BERT October 11, 2018	RoBERTa July 26, 2019	DistilBERT October 2, 2019	ALBERT September 26, 2019
Parameters	Base: 110M Large: 340M	Base: 125 Large: 355	Base: 66	Base: 12M Large: 18M
Layers / Hidden Dimensions / Self-Attention Heads	Base: 12 / 768 / 12 Large: 24 / 1024 / 16	Base: 12 / 768 / 12 Large: 24 / 1024 / 16	Base: 6 / 768 / 12	Base: 12 / 768 / 12 Large: 24 / 1024 / 16
Training Time	Base: 8 x V100 x 12d Large: 280 x V100 x 1d	1024 x V100 x 1 day (4-5x more than BERT)	Base: 8 x V100 x 3.5d (4 times less than BERT)	[not given]  Large: 1.7x faster
Performance	Outperforming SOTA in Oct 2018	88.5 on GLUE	97% of BERT-base's performance on GLUE	89.4 on GLUE
Pre-Training Data	BooksCorpus + English Wikipedia = 16 GB	BERT + CCNews + OpenWebText + Stories = 160 GB	BooksCorpus + English Wikipedia = 16 GB	BooksCorpus + English Wikipedia = 16 GB
Method	Bidirectional Transformer, MLM & NSP	BERT without NSP, Using Dynamic Masking	BERT Distillation	BERT with reduced parameters & SOP (not NSP)

### 3 BERT

We introduce BERT and its detailed implementation in this section. There are two steps in our framework: *pre-training* and *fine-tuning*. During pre-training, the model is trained on unlabeled data over different pre-training tasks. For fine-tuning, the BERT model is first initialized with the pre-trained parameters, and all of the parameters are fine-tuned using labeled data from the downstream tasks. Each downstream task has separate fine-tuned models, even though they are initialized with the same pre-trained parameters. The question-answering example in Figure 1 will serve as a running example for this section.

A distinctive feature of BERT is its unified architecture across different tasks. There is mini-

<sup>1</sup>In this work, we denote the number of layers (i.e., Transformer blocks) as  $L$ , the hidden size as  $H$ , and the number of self-attention heads as  $A$ .<sup>2</sup> We primarily report results on two model sizes: **BERT<sub>BASE</sub>** ( $L=12$ ,  $H=768$ ,  $A=12$ , Total Parameters=110M) and **BERT<sub>LARGE</sub>** ( $L=24$ ,  $H=1024$ ,  $A=16$ , Total Parameters=340M).

BERT<sub>BASE</sub> was chosen to have the same model size as OpenAI GPT for comparison purposes. Critically, however, the BERT Transformer uses bidirectional self-attention, while the GPT Transformer uses constrained self-attention where every token can only attend to context to its left.<sup>4</sup>

<sup>1</sup><https://github.com/tensorflow/tensor2tensor>

<sup>2</sup><http://nlp.seas.harvard.edu/2018/04/03/attention.html>

<sup>3</sup>In all cases we set the feed-forward/filter size to be  $4H$ , i.e., 3072 for the  $H = 768$  and 4096 for the  $H = 1024$ .

<sup>4</sup>We note that in the literature the bidirectional Trans-

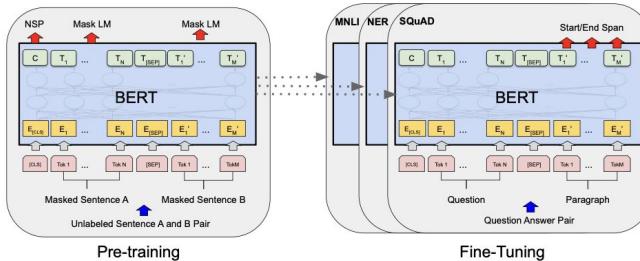


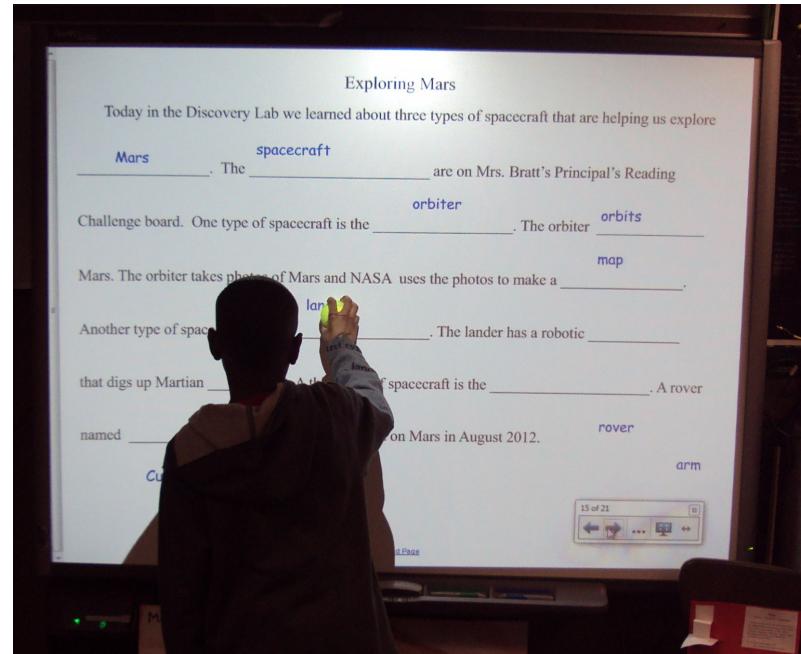
Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different downstream tasks. During fine-tuning, all parameters are fine-tuned. [CITATION] is a special



# A cloze test (Taylor, 1953) and Masked Language Model (MLM)

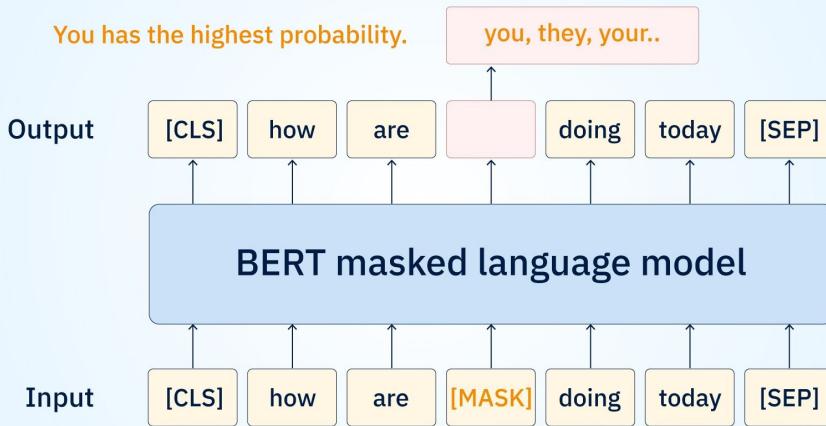
**Input/Output Representations** To make BERT handle a variety of down-stream tasks, our input representation is able to unambiguously represent both a single sentence and a pair of sentences (e.g., (Question, Answer)) in one token sequence. Throughout this work, a “sentence” can be an arbitrary span of contiguous text, rather than an actual linguistic sentence. A “sequence” refers to the input token sequence to BERT, which may be a single sentence or two sentences packed together.

In order to train a deep bidirectional representation, we simply mask some percentage of the input tokens at random, and then predict those masked tokens. We refer to this procedure as a “masked LM” (MLM), although it is often referred to as a Cloze task in the literature (Taylor, 1953). In this case, the final hidden vectors corresponding to the mask tokens are fed into an output softmax over the vocabulary, as in a standard LM. In all of our experiments, we mask 15% of all WordPiece to-



A cloze test is an exercise, test, or assessment in which a portion of text is masked and the participant is asked to fill in the masked portion of text . Source: wikipedia

# Masked Language Model (MLM)



Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	#ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{\text{my}}$	$E_{\text{dog}}$	$E_{\text{is}}$	$E_{\text{cute}}$	$E_{[\text{SEP}]}$	$E_{\text{he}}$	$E_{\text{likes}}$	$E_{\text{play}}$	$E_{\#ing}$	$E_{[\text{SEP}]}$
Segment Embedments	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

Figure 2: BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.

The NSP task is closely related to representation-learning objectives used in Jernite et al. (2017) and Logeswaran and Lee (2018). However, in prior work, only sentence embeddings are transferred to downstream tasks, such as sequence tagging or question answering, and the  $[\text{CLS}]$  representation is fed into an output layer for classification, such as entailment or sentiment analysis.

**Pre-training data** The pre-training procedure largely follows the existing literature on language model pre-training. For the pre-training corpus we use the BooksCorpus (800M words) (Zhu et al., 2015) and English Wikipedia (2,500M words). For Wikipedia we extract only the text passages and ignore lists, tables, and headers. It is critical to use a document-level corpus rather than a shuffled sentence-level corpus such as the Billion Word Benchmark (Chelba et al., 2013) in order to extract long contiguous sequences.

## 3 Experiments

In this section, we present BERT fine-tuning results on 11 NLP tasks.

### 3.2 Fine-tuning BERT

Fine-tuning is straightforward since the self-attention mechanism in the Transformer allows BERT to model many downstream tasks—whether they involve single text or text pairs—by swapping out the appropriate inputs and outputs. For applications involving text pairs, a common pattern is to independently encode text pairs before applying bidirectional cross attention, such as Parikh et al. (2016); Seo et al. (2017). BERT instead uses the self-attention mechanism to unify these two stages, as encoding a concatenated text pair with self-attention effectively includes *bidirectional* cross attention between two sentences.

For each task, we simply plug in the task-specific inputs and outputs into BERT and fine-tune all the parameters end-to-end. At the input, sentence A and sentence B from pre-training are analogous to (1) sentence pairs in paraphrasing, (2) hypothesis-premise pairs in entailment, (3) question-passage pairs in question answering, and

(4) a degenerate text-∅ pair in text classification or sequence tagging. At the output, the token representations are fed into an output layer for token-level tasks, such as sequence tagging or question answering, and the  $[\text{CLS}]$  representation is fed into an output layer for classification, such as entailment or sentiment analysis.

Compared to pre-training, fine-tuning is relatively inexpensive. All of the results in the paper can be replicated in at most 1 hour on a single Cloud TPU, or a few hours on a GPU, starting from the exact same pre-trained model.<sup>7</sup> We describe the task-specific details in the corresponding subsections of Section 4. More details can be found in Appendix A.5.

## 4 Experiments

In this section, we present BERT fine-tuning results on 11 NLP tasks.

### 4.1 GLUE

The General Language Understanding Evaluation (GLUE) benchmark (Wang et al., 2018a) is a collection of diverse natural language understanding tasks. Detailed descriptions of GLUE datasets are included in Appendix B.1.

To fine-tune on GLUE, we represent the input sequence (for single sentence or sentence pairs) as described in Section 3, and use the final hidden vector  $C \in \mathbb{R}^H$  corresponding to the first input token ( $[\text{CLS}]$ ) as the aggregate representation. The only new parameters introduced during fine-tuning are classification layer weights  $W \in \mathbb{R}^{K \times H}$ , where  $K$  is the number of labels. We compute a standard classification loss with  $C$  and  $W$ , i.e.,  $\log(\text{softmax}(CW^T))$ .

<sup>7</sup>For example, the BERT SQuAD model can be trained in around 30 minutes on a single Cloud TPU to achieve a Dev F1 score of 91.0%.

<sup>8</sup>See (10) in <https://gluebenchmark.com/faq>.

## Pre-training objectives

- Masked Language Modeling (MLM): mask ~15% of tokens, then predict the originals using both left and right context.
- Next Sentence Prediction (NSP, BERT): classify whether Sentence B truly follows Sentence A (IsNext vs NotNext).
- Scale matters: train on a large, diverse text corpus to capture broad linguistic patterns.

## What the model sees

### MLM EXAMPLE

Original: The quick brown fox jumps over the lazy dog.

Input: The quick brown [MASK] jumps over the lazy dog.

Predict: "fox"

### NSP EXAMPLE (BERT)

Sentence A: The quick brown fox jumps over the lazy dog.

Sentence B: It then swiftly ran up the hill.

Classifier: IsNext / NotNext

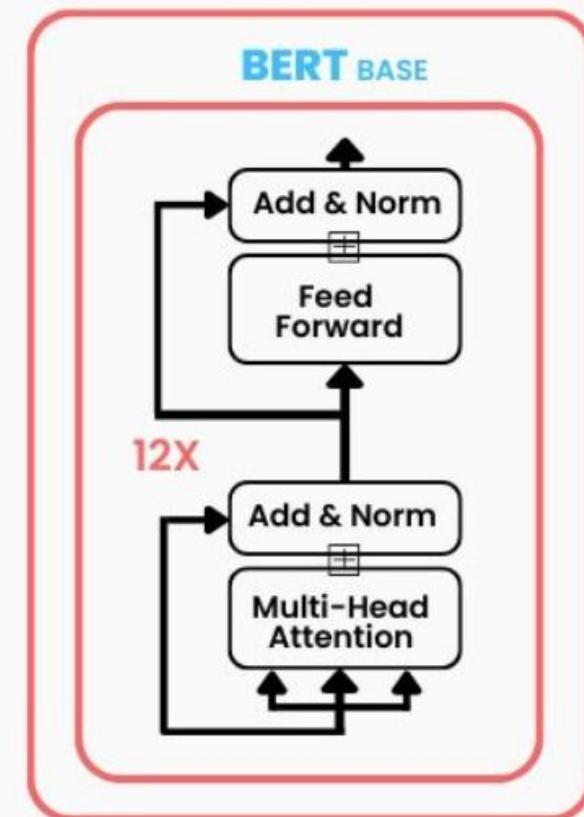
```

BertModel(
    (embeddings): BertEmbeddings(
        (word_embeddings): Embedding(30522, 768, padding_idx=0)
        (position_embeddings): Embedding(512, 768)
        (token_type_embeddings): Embedding(2, 768)
        (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
        (dropout): Dropout(p=0.1, inplace=False)
    )
    (encoder): BertEncoder(
        (layer): ModuleList(
            (0-11): 12 x BertLayer(
                (attention): BertAttention(
                    (self): BertSelfAttention(
                        (query): Linear(in_features=768, out_features=768, bias=True)
                        (key): Linear(in_features=768, out_features=768, bias=True)
                        (value): Linear(in_features=768, out_features=768, bias=True)
                        (dropout): Dropout(p=0.1, inplace=False)
                    )
                )
                (output): BertSelfOutput(
                    (dense): Linear(in_features=768, out_features=768, bias=True)
                    (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
                    (dropout): Dropout(p=0.1, inplace=False)
                )
            )
            (intermediate): BertIntermediate(
                (dense): Linear(in_features=768, out_features=3072, bias=True)
                (intermediate_act_fn): GELUActivation()
            )
            (output): BertOutput(
                (dense): Linear(in_features=3072, out_features=768, bias=True)
                (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
                (dropout): Dropout(p=0.1, inplace=False)
            )
        )
    )
)

```

## Lecture 2

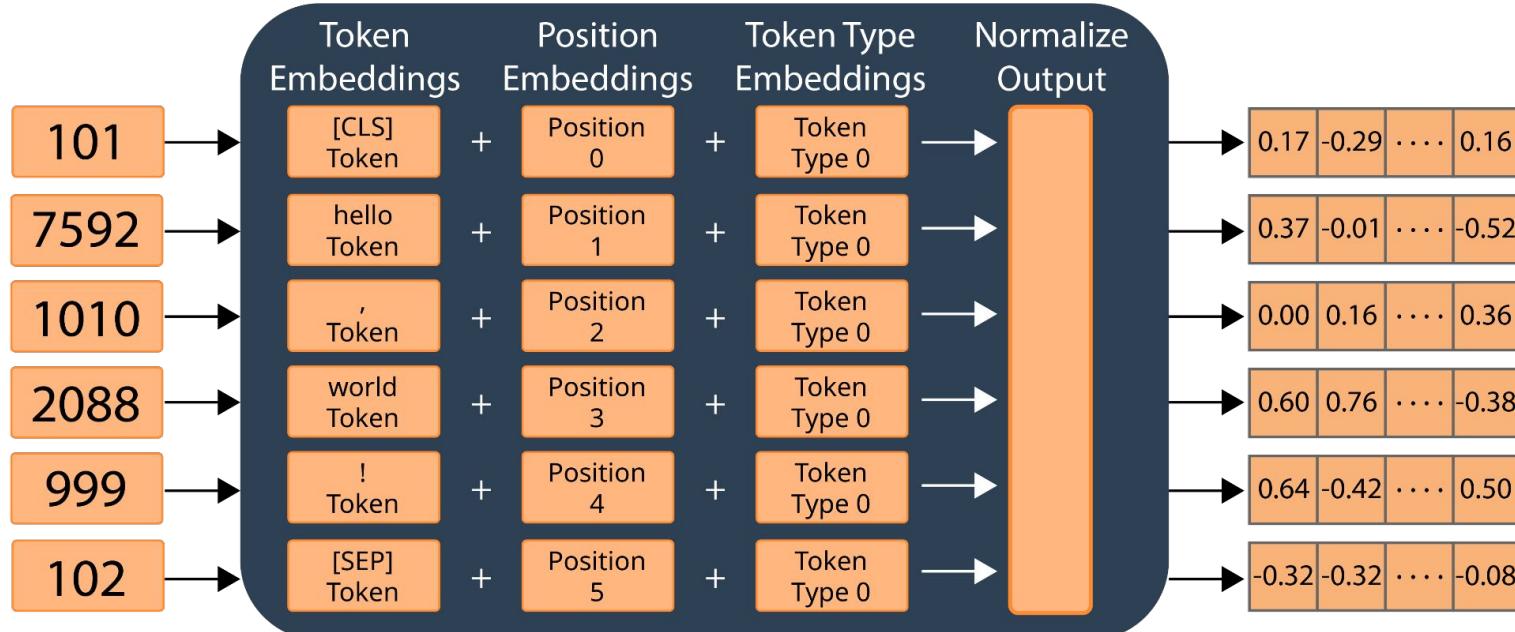
### Lecture 1



110M Parameters

# BERT Embedding Layer

BERT Embedding Layer



# Absolute vs Relative Positioning

Sequence	Index of token, $k$	Positional Encoding Matrix with $d=4$ , $n=100$			
		$i=0$	$i=0$	$i=1$	$i=1$
I	0	$P_{00}=\sin(0) = 0$	$P_{01}=\cos(0) = 1$	$P_{02}=\sin(0) = 0$	$P_{03}=\cos(0) = 1$
am	1	$P_{10}=\sin(1/1) = 0.84$	$P_{11}=\cos(1/1) = 0.54$	$P_{12}=\sin(1/10) = 0.10$	$P_{13}=\cos(1/10) = 1.0$
a	2	$P_{20}=\sin(2/1) = 0.91$	$P_{21}=\cos(2/1) = -0.42$	$P_{22}=\sin(2/10) = 0.20$	$P_{23}=\cos(2/10) = 0.98$
Robot	3	$P_{30}=\sin(3/1) = 0.14$	$P_{31}=\cos(3/1) = -0.99$	$P_{32}=\sin(3/10) = 0.30$	$P_{33}=\cos(3/10) = 0.96$

Positional Encoding Matrix for the sequence 'I am a robot'

# Absolute vs Relative Positioning

## Modeling Graph Structure via Relative Position for Text Generation from Knowledge Graphs

Martin Schmitt<sup>1</sup> Leonardo F. R. Ribeiro<sup>2</sup> Philipp Dufter<sup>1</sup> Iryna Gurevych<sup>2</sup> Hinrich Schütze<sup>1</sup>

<sup>1</sup>Center for Information and Language Processing (CIS), LMU Munich

<sup>2</sup>Research Training Group AIPHES and UKP Lab, Technische Universität Darmstadt

[martin@cis.lmu.de](mailto:martin@cis.lmu.de)

## Abstract

et al., 2018) or variants of Transformer (Vaswani et al., 2017) that apply self-attention on all nodes together, including those that are not connected. To avoid losing information, approaches use

Published at the MLDD workshop, ICLR 2021

# Do Transformers Really Perform Bad for Graph Representation?

# GRPE: RELATIVE POSITIONAL ENCODING FOR GRAPH TRANSFORMER

\* Chang\*, <sup>1</sup>Donggeon Lee, <sup>1</sup>Juntae Kim, <sup>2</sup>Seung-won Hwang

## ABSTRACT

**ABSTRACT**

Model to encode graphs is a key challenge of molecular processing, but it requires explicit incorporation of position bias terms. Existing approaches either linearize a graph to encode atom in the sequence of nodes, or encode relative position with angularization, while the latter loses a tight integration of node-edge and graph information. In this work, we propose relative positional encoding to overcome the weakness of the previous approaches. Our method codes a graph without linearization and considers both node-spatial relation and node-edge relation. We name our method Graph Relative Positional Encoding dedicated to graph representation learning. Experiments conducted on various molecular property prediction datasets show that the proposed method outperforms previous approaches significantly. Our code is publicly available at <https://github.com/lenscloth/GRPE>.

06.05234v5 [cs.LG] 24 Nov 2021

# Socioeconomic phenomena as sequential processes

Many human phenomena unfold as ordered sequences rather than isolated events  
The meaning of an action often depends on what comes before and after it

## Examples

Language: Word order shapes meaning and grammar. The same words in a different sequence can change or destroy meaning

Grandma, let's eat! ☐ Inviting grandma to dinner.

Let's eat grandma! ☐ Cannibalism 😱

More examples?

# Socioeconomic phenomena as sequential processes

Many human phenomena unfold as ordered sequences rather than isolated events  
The meaning of an action often depends on what comes before and after it

## Examples

Language: Word order shapes meaning and grammar. The same words in a different sequence can change or destroy meaning

## Life course events

Typical sequences such as kindergarten, primary school, secondary school, and so on.

## Firms and products

Products are developed, launched, and upgraded in specific sequences. Earlier technological choices constrain or enable future innovations

# Life2Vec

nature computational science

Article

<https://doi.org/10.1038/s43588-023-00573-5>

## Using sequences of life-events to predict human lives

Received: 6 June 2023

Accepted: 15 November 2023

Published online: 18 December 2023

 Check for updates

Germane Sovisens<sup>Ø</sup><sup>1</sup>, Tina Eliassi-Rad<sup>Ø,2,3</sup>, Lars Kai Hansen<sup>1</sup>, Laust Hvas Mortensen<sup>Ø,4,5</sup>, Lau Lilleholt<sup>Ø,6,7</sup>, Anna Rogers<sup>8</sup>, Ingo Zettler<sup>Ø,6,7</sup> & Sune Lehmann<sup>Ø,1,7,9</sup>

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

We live in the age of algorithm-driven prediction of human behavior. The predictions range from those at the global and population level, with societies allocating vast resources to predicting phenomena such as global warming or the spread of infectious disease, all the way to the commercialization of individual predictions that tell us what to eat and how to behave as we use social media. When it comes to individual life outcomes, however, the picture is more complex. Sociodemographic factors play an important role in human lives<sup>1–3</sup>, but, based on independent analyses of the same dataset, a recent collaboration of 160 teams has recently argued for practical upper limits for the predictions of life outcomes<sup>4</sup>.

In this Article we find that, with highly detailed records, a different picture of individual-level predictability emerges. Drawing on a unique dataset consisting of detailed individual-level day-by-day records<sup>5–7</sup> describing the six million inhabitants of Denmark, and spanning a

decade interval, we show that accurate individual predictions are indeed possible. Our dataset includes a host of indicators, such as health, professional occupation and affiliation, income level, residency, working hours and education, and is used to predict a range of outcomes of individual lives.

The challenge here is that we are not only experiencing this “age of human prediction” as the advent of massive datasets and powerful machine learning algorithms<sup>8–10</sup>. Over the past decade, machine learning has revolutionized the image- and text-processing fields by accessing ever larger datasets that have enabled increasingly complex models<sup>11–13</sup>. Language processing has evolved particularly rapidly, and transformer architectures have proven successful at capturing complex patterns in language and structured knowledge<sup>14–16</sup>. These language models optimized in natural language processing, the ability to capture structure in human language generalize to other sequences<sup>17–19</sup>, that share properties with language, for example, where sequence

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

<sup>1</sup>Nano Computer, Technical University of Denmark, Lyngby, Denmark, <sup>2</sup>Network Science Institute, Northeastern University, Boston, MA, USA,

<sup>3</sup>Faculty College of Computer Sciences, Northeastern University, Boston, MA, USA, <sup>4</sup>Data Science Lab, Statistic Denmark, Copenhagen, Denmark,

<sup>5</sup>Department of Public Health, University of Copenhagen, Copenhagen, Denmark, <sup>6</sup>Department of Psychology, University of Copenhagen, Copenhagen, Denmark, <sup>7</sup>Copenhagen Center for Social Data Science (CCSDS), University of Copenhagen, Copenhagen, Denmark, <sup>8</sup>Computer Science Department,

<sup>9</sup>TU University of Copenhagen, Copenhagen, Denmark, <sup>10</sup>e-mail: stp@dtu.dk

# Life2Vec

nature computational science

Article

<https://doi.org/10.1038/s43588-023-00573-5>

## Using sequences of life-events to predict human lives

Received: 6 June 2023

Accepted: 15 November 2023

Published online: 18 December 2023

 Check for updates

Germane Sovisens<sup>Ø</sup><sup>1</sup>, Tina Eliassi-Rad<sup>Ø,2,3</sup>, Lars Kai Hansen<sup>1</sup>, Laust Hvas Mortensen<sup>Ø,4,5</sup>, Lau Lilleholt<sup>Ø,6,7</sup>, Anna Rogers<sup>8</sup>, Ingo Zettler<sup>Ø,6,7</sup> & Sune Lehmann<sup>Ø,1,7,9</sup>

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

We live in the age of algorithm-driven prediction of human behavior. The predictions range from those at the global and population level, with societies allocating vast resources to predicting phenomena such as global warming or the spread of infectious disease, all the way to the commercialization of individual predictions that tell us what to eat and how to behave as we use social media. When it comes to individual life outcomes, however, the picture is more complex. Sociodemographic factors play an important role in human lives<sup>1–3</sup>, but, based on independent analyses of the same dataset, a recent collaboration of 160 teams has recently argued for practical upper limits for the predictions of life outcomes<sup>4</sup>.

In this Article we find that, with highly detailed records, a different picture of individual-level predictability emerges. Drawing on a unique dataset consisting of detailed individual-level day-by-day records<sup>5–7</sup> describing the six million inhabitants of Denmark, and spanning a

decade interval, we show that accurate individual predictions are indeed possible. Our dataset includes a host of indicators, such as health, professional occupation and affiliation, income level, residency, working hours and education, and is used here to predict human behavior as we use social media.

When it comes to individual life outcomes, however, we are currently experiencing this “age of human prediction” as the advent of massive datasets and powerful machine learning algorithms<sup>8–10</sup>. Over the past decade, machine learning has revolutionized the image- and text-processing fields by accessing ever larger datasets that have enabled increasingly complex models<sup>11,12</sup>. Language processing has evolved particularly rapidly, and transformer architectures have proven successful at capturing complex patterns in language and structured knowledge<sup>13–15</sup>. These language models outperform traditional language processing in their ability to capture structure in human language generalizes to other sequences<sup>16–19</sup>, that share properties with language, for example, where sequence

<sup>1</sup>Nano Computer, Technical University of Denmark, Lyngby, Denmark, <sup>2</sup>Network Science Institute, Northeastern University, Boston, MA, USA,

<sup>3</sup>Brookings College of Computer Sciences, Northeastern University, Boston, MA, USA, <sup>4</sup>Data Science Lab, Statistic Denmark, Copenhagen, Denmark,

<sup>5</sup>Department of Public Health, University of Copenhagen, Copenhagen, Denmark, <sup>6</sup>Department of Psychology, University of Copenhagen, Copenhagen, Denmark, <sup>7</sup>Copenhagen Center for Social Data Science (CCSDS), University of Copenhagen, Copenhagen, Denmark, <sup>8</sup>Computer Science Department, IT University of Copenhagen, Copenhagen, Denmark, <sup>9</sup>e-mail: stø@itu.dk

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

# Life2Vec

nature computational science

Article

<https://doi.org/10.1038/s43588-023-00573-5>

## Using sequences of life-events to predict human lives

Received: 6 June 2023

Accepted: 15 November 2023

Published online: 18 December 2023

 Check for updates

Germane Sovisens<sup>Ø</sup><sup>1</sup>, Tina Eliassi-Rad<sup>Ø,2,3</sup>, Lars Kai Hansen<sup>1</sup>, Laust Hvas Mortensen<sup>Ø,4,5</sup>, Lau Lilleholt<sup>Ø,6,7</sup>, Anna Rogers<sup>8</sup>, Ingo Zettler<sup>Ø,6,7</sup> & Sune Lehmann<sup>Ø,1,7</sup>

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

We live in the age of algorithm-driven prediction of human behavior. The predictions range from those at the global and population level, with societies allocating vast resources to predicting phenomena such as global warming or the spread of infectious diseases, all the way to the computing that informs individual predictions that tell us what to eat and how to act as we use social media. When it comes to individual life outcomes, however, the picture is more complex. Sociodemographic factors play an important role in human lives<sup>1–3</sup>, but, based on independent analyses of the same dataset, a recent collaboration of 160 teams has recently argued for practical upper limits for the predictions of life outcomes<sup>4</sup>.

In this Article we find that, with highly detailed data, a different picture of individual-level predictability emerges. Drawing on a unique dataset consisting of detailed individual-level day-by-day records<sup>5–7</sup> describing the six million inhabitants of Denmark, and spanning a

decade interval, we show that accurate individual predictions are indeed possible. Our dataset includes a host of indicators, such as health, professional occupation and affiliation, income level, residency, working hours and education, and is based on administrative records.

The question now is how we are currently experiencing this “age of human prediction”: is the advent of massive datasets and powerful machine learning algorithms<sup>8–10</sup>? Over the past decade, machine learning has revolutionized the image- and text-processing fields by accessing ever larger datasets that have enabled increasingly complex models<sup>11,12</sup>. Language processing has evolved particularly rapidly, and transformer architectures have proven successful at capturing complex patterns in text and structured language<sup>13–16</sup>. A key feature of these models is that they are trained on large amounts of text that share properties with language, for example, where sequence

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

<sup>1</sup>Nano Computer, Technical University of Denmark, Lyngby, Denmark, <sup>2</sup>Network Science Institute, Northeastern University, Boston, MA, USA,

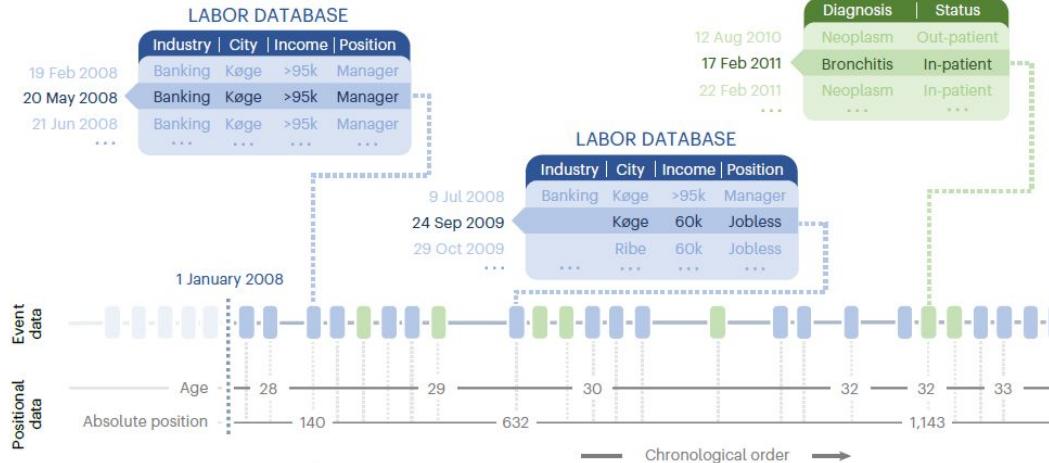
<sup>3</sup>Faculty College of Computer Sciences, Northeastern University, Boston, MA, USA, <sup>4</sup>Data Science Lab, Statistic Denmark, Copenhagen, Denmark,

<sup>5</sup>Department of Public Health, University of Copenhagen, Copenhagen, Denmark, <sup>6</sup>Department of Psychology, University of Copenhagen, Copenhagen, Denmark, <sup>7</sup>Copenhagen Center for Social Data Science (CCSDS), University of Copenhagen, Copenhagen, Denmark, <sup>8</sup>Computer Science Department,

<sup>9</sup>TU University of Copenhagen, Copenhagen, Denmark, <sup>10</sup>e-mail: stj@dtu.dk

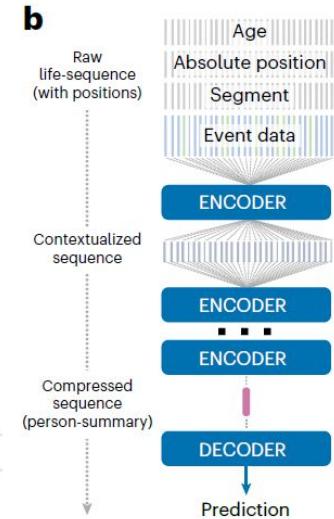
# Life2Vec

**a**



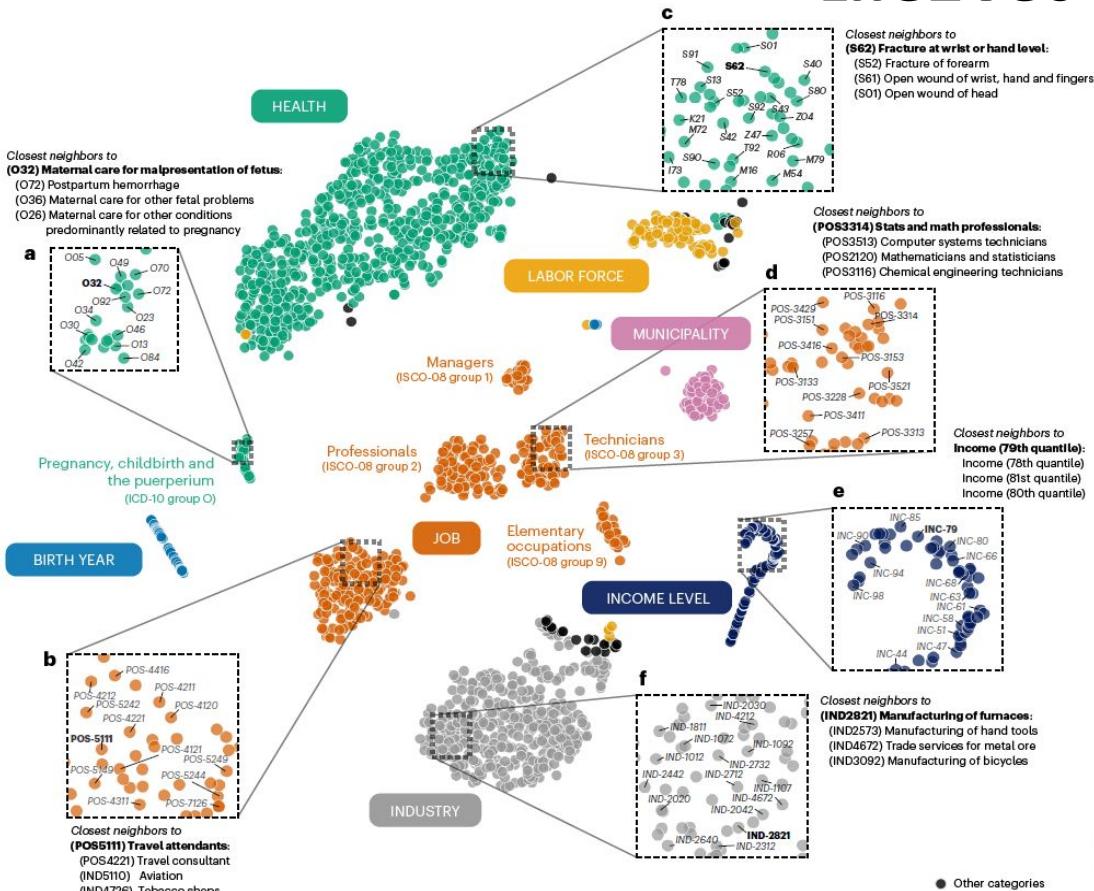
**Fig. 1 | A schematic individual-level data representation for the life2vec model.** **a,b**, We organize socio-economic and health data from the Danish national registers from 1 January 2008 to 31 December 2015 into a single chronologically ordered life-sequence (**a**). Each database entry becomes an event in the sequence, where an event has associated positional and contextual data. The contextual data include variables associated with the entry (for example, industry, city, income and job type). The positional data include the person's age (expressed in full years) and absolute position (number of days since 1

**b**



January 2008). The raw life-sequence is then passed to the model described in **b**. The model consists of multiple stacked encoders. The first encoder combines contextual and positional information to produce a contextual representation of each life-event. The following encoders output deep contextual representations of each life-event (considering the overall content of the life-sequence). The final encoder layer fuses the representations of life-events to produce the representation of a life-sequence. The decoder uses the latter to make predictions.

# Life2Vec



**Fig. 2 | Two-dimensional projection of the concept space (using PaCMAP).** Each point corresponds to a concept token in the vocabulary ( $n = 2,043$ ). Points are colored based on the concept types (infrequent types are represented as black points). Each region provides a zoom of a part of the concept space. The top three closest neighbors for selected tokens (based on the cosine distance) are also displayed. **a**, Diagnoses related to pregnancy, childbirth and the puerperium

**b**, Job concepts related to service and sales workers (corresponds to Job category 5 of ISCO-08<sup>25</sup>). **c**, Injury-related diagnoses in ICD-10<sup>27</sup>. **d**, Job concepts related to technicians and associate professionals (corresponds to Job category 3 of ISCO-08<sup>25</sup>). **e**, Income-related concepts. Life2vec arranges these concepts in increasing ordinal order. **f**, Concepts related to the manufacturing industry in DB07<sup>26</sup>.

# Life2Vec

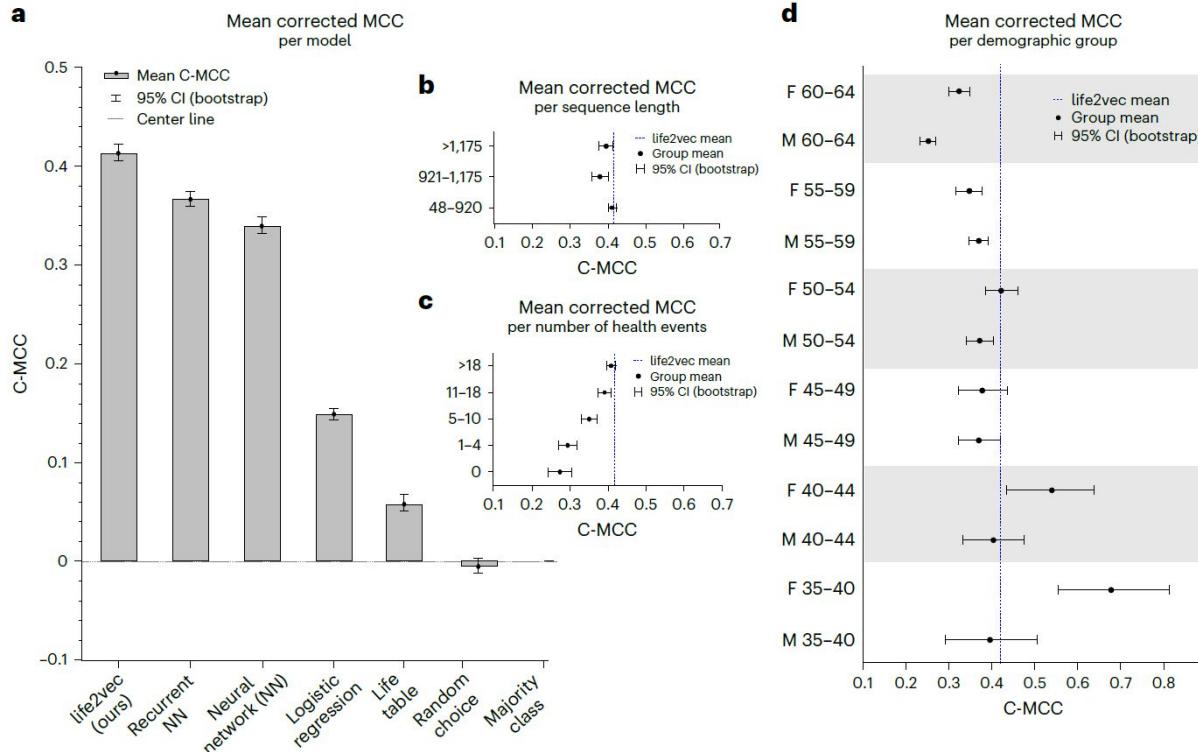
## Matthews Correlation Coefficient (MCC)

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Where:

- **TP** = True Positives
  - **TN** = True Negatives
  - **FP** = False Positives
  - **FN** = False Negatives
- 
- MCC considers all four outcomes of the confusion matrix
  - Values range from  $-1$  to  $+1$
  - $+1$  indicates perfect prediction
  - $0$  indicates random prediction
  - $-1$  indicates total disagreement

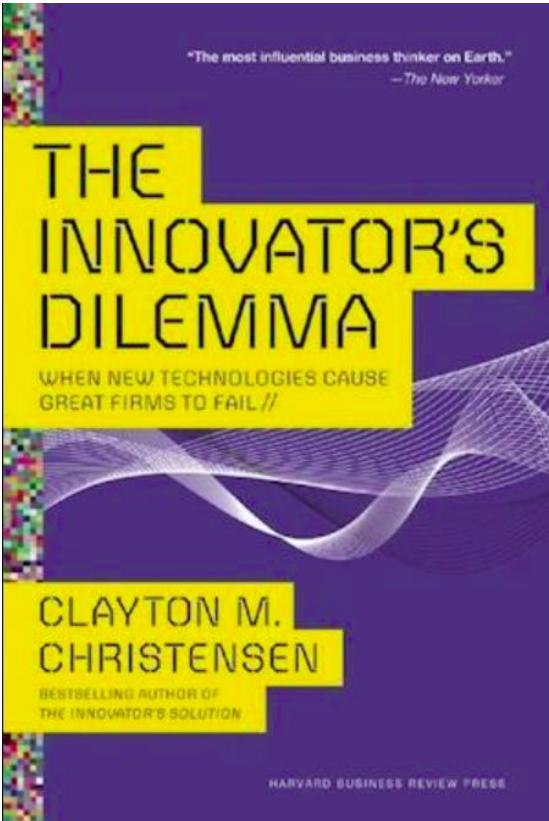
# Life2Vec



**Fig. 3 | Performance of models on the mortality prediction task quantified with the mean C-MCC with 95% confidence interval.** **a**, Comparison of life2vec performance to baselines ( $n = 100,000$ ). **b-d**, Performance of life2vec on different cohorts of the population: performance of life2vec per sequence length

(**b**), performance of life2vec based on the number of health events in a sequence (**c**) and performance of life2vec per intersectional group (based on age group and sex) (**d**). F, female; M, male.

# Transformers and BDS



Successful firms often fail:

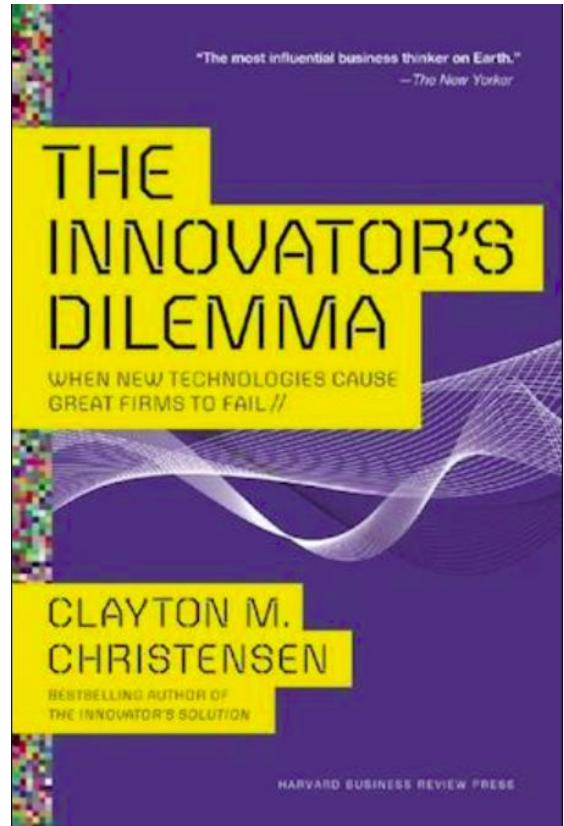
- not because they are poorly managed
- because they focus too strongly on existing customers and current technologies

Disruptive innovations start as simpler, cheaper, or lower-performance alternatives

- Disruptive innovations initially serve niche or low-end markets
- Over time, they improve and displace established technologies and firms
- The book conceptualizes the role of timing, sequencing, and path dependency in innovation

# Transformers and BDS

## Related and unrelated diversification (remember?)



### The Principle of Relatedness

César A. Hidalgo<sup>1</sup><sup>(✉)</sup>, Pierre-Alexandre Balland<sup>2</sup>, Ron Boschma<sup>2,3</sup>,  
Mercedes Delgado<sup>4</sup>, Maryann Feldman<sup>5</sup>, Koen Frenken<sup>6</sup>,  
Edward Glaeser<sup>7,8</sup>, Canfei He<sup>9</sup>, Dieter F. Kogler<sup>10</sup>, Andrea Morrison<sup>2</sup>,  
Frank Nettek<sup>11</sup>, David Rigby<sup>12</sup>, Scott Stern<sup>4,8</sup>, Siqui Zheng<sup>13,14</sup>,  
and Shengjun Zhu<sup>9</sup>

<sup>1</sup> Collective Learning Group, the MIT Media Lab,  
Massachusetts Institute of Technology, Cambridge, USA  
hidalgo@mit.edu

<sup>2</sup> Department of Human Geography and Planning,  
Utrecht University, Utrecht, Netherlands

<sup>3</sup> Business School,

University of Stavanger, Stavanger, Norway

<sup>4</sup> MIT Sloan School of Management, Cambridge, USA

<sup>5</sup> Department of Public Policy,

University of North Carolina at Chapel Hill, Chapel Hill, USA

<sup>6</sup> Copernicus Institute of Sustainable Development,

Utrecht University, Utrecht, Netherlands

<sup>7</sup> Department of Economics, Harvard University, Cambridge, USA

<sup>8</sup> National Bureau of Economic Research, Cambridge, USA

<sup>9</sup> College of Urban and Environmental Sciences,

Peking University, Beijing, China

<sup>10</sup> Spatial Dynamics Lab and School of Architecture Planning  
and Environmental Policy, University College Dublin, Dublin, Ireland

<sup>11</sup> Center for International Development, Harvard University, Cambridge, USA

<sup>12</sup> Department of Geography and Department of Statistics,

University of California, Los Angeles, USA

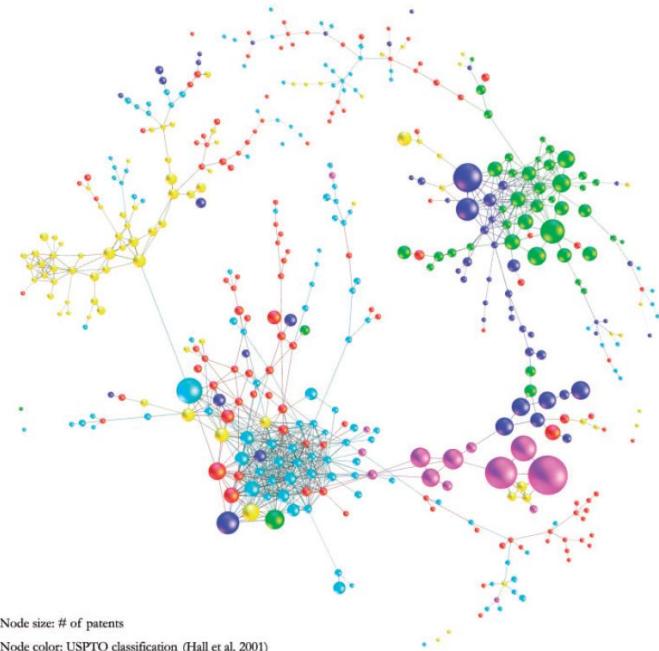
<sup>13</sup> Department of Urban Studies and Planning, MIT, Cambridge, USA

<sup>14</sup> School of Civil Engineering, Tsinghua University, Beijing, China

**Abstract.** The idea that skills, technology, and knowledge, are spatially concentrated, has a long academic tradition. Yet, only recently this hypothesis has been empirically formalized and corroborated at multiple spatial scales, for different economic activities, and for a diversity of institutional regimes. The new synthesis is an empirical principle describing the probability that a region enters—or exits—an economic activity as a function of the number of related activities present in that location. In this paper we summarize some of the recent empirical evidence that has generalized the principle of relatedness to a fact describing the entry and exit of products, industries, occupations, and technologies, at the national, regional, and metropolitan scales. We conclude by describing some of the policy implications and future avenues of research implied by this robust empirical principle.

C. A. Hidalgo and P.-A. Balland—Contributed equally.

© Springer Nature Switzerland AG 2018  
A. J. Morales et al. (Eds.): ICCS 2018, SPCM, pp. 451–457, 2018.  
[https://doi.org/10.1007/978-3-319-96661-8\\_46](https://doi.org/10.1007/978-3-319-96661-8_46)



Node size: # of patents

Node color: USPTO classification (Hall et al. 2001)

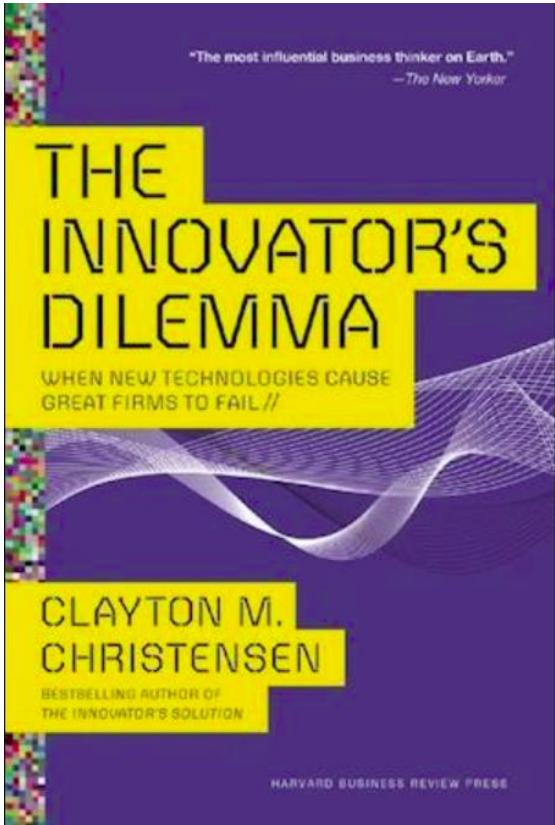
Mechanical      Electrical and Electronic

Chemical      Computers and Communications

Drugs and Medical      Others

# Transformers and BDS

## Related and unrelated diversification (remember?)



### The Principle of Relatedness

César A. Hidalgo<sup>1,✉</sup>, Pierre-Alexandre Balland<sup>2</sup>, Ron Boschma<sup>2,3</sup>,  
Mercedes Delgado<sup>4</sup>, Maryann Feldman<sup>5</sup>, Koen Frenken<sup>6</sup>,  
Edward Glaeser<sup>7,8</sup>, Canfei He<sup>9</sup>, Dieter F. Kogler<sup>10</sup>, Andrea Morrison<sup>2</sup>,  
Frank Nettek<sup>11</sup>, David Rigby<sup>12</sup>, Scott Stern<sup>4,8</sup>, Siqui Zheng<sup>13,14</sup>,  
and Shengjun Zhu<sup>9</sup>

<sup>1</sup> Collective Learning Group, the MIT Media Lab,  
Massachusetts Institute of Technology, Cambridge, USA  
✉ [catalin.hidalgo@mit.edu](mailto:catalin.hidalgo@mit.edu)

relatedness among activities. They look, for instance, at the co-export of products [2], the flow of labor among industries [3], or combined measures of input-output links and shared labor pools [4]. This methodological flexibility, has allowed scholars from a variety of fields to document a robust and reproducible relationship between the probability that a location will develop expertise in a new industry [5–7], technology [8, 9], research area [10], product [2], or occupation [11], and the number of related activities that are already present in that location.

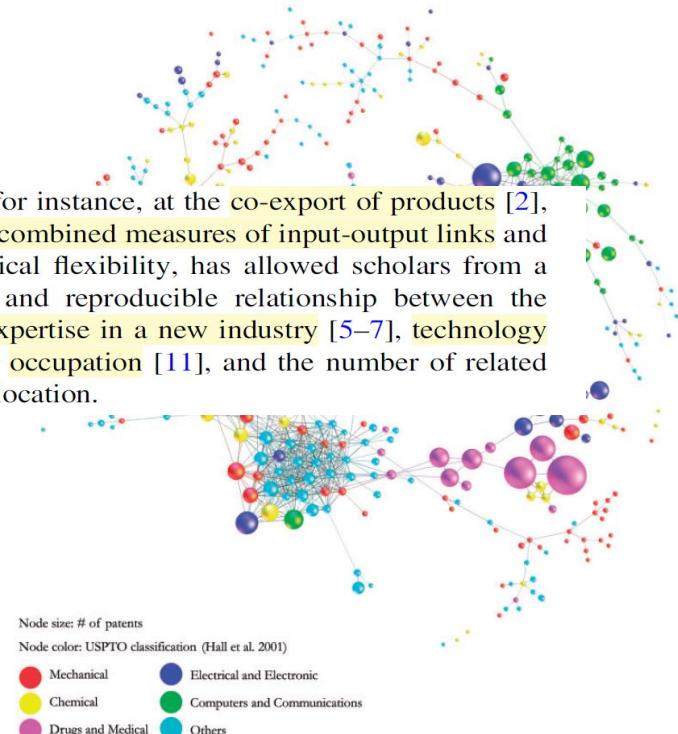
Hidalgo et al. 2018 (p. 452)

**Abstract.** The idea that skills, technology, and knowledge, are spatially concentrated, has a long academic tradition. Yet, only recently this hypothesis has been empirically formalized and corroborated at multiple spatial scales, for different economic activities, and for a diversity of institutional regimes. The new synthesis is an empirical principle describing the probability that a region enters—or exits—an economic activity as a function of the number of related activities present in that location. In this paper we summarize some of the recent empirical evidence that has generalized the principle of relatedness to a fact describing the entry and exit of products, industries, occupations, and technologies, at the national, regional, and metropolitan scales. We conclude by describing some of the policy implications and future avenues of research implied by this robust empirical principle.

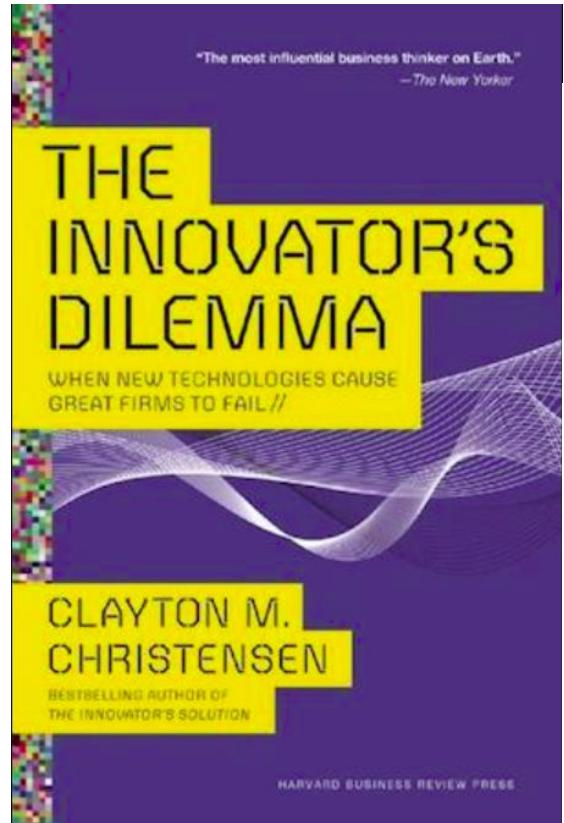
C. A. Hidalgo and P.-A. Balland—Contributed equally.

© Springer Nature Switzerland AG 2018

A. J. Morales et al. (Eds.): ICCS 2018, SPCM, pp. 451–457, 2018.  
[https://doi.org/10.1007/978-3-319-96661-8\\_46](https://doi.org/10.1007/978-3-319-96661-8_46)



# Transformers and BDS



## Related and unrelated diversification (remember?)



### Related Variety and Regional Development: A Critique

**Harald Barthelt** Department of Geography and Planning University of Toronto Toronto, Ontario M5S 3G3 Canada harald.barthelt@utoronto.ca

**Michael Storper** Department of Geography and Environment London School of Economics London WC2A 2AE UK and London School of Public Affairs UCLA Los Angeles, CA 90095-1656 m.storper@lse.ac.uk

**Key words:**  
economic  
geographies  
of places  
evolutionary  
economic  
geography (EEG)  
regional development  
regional specialization  
related variety

**JEL codes:**  
L23  
R11

#### abstract

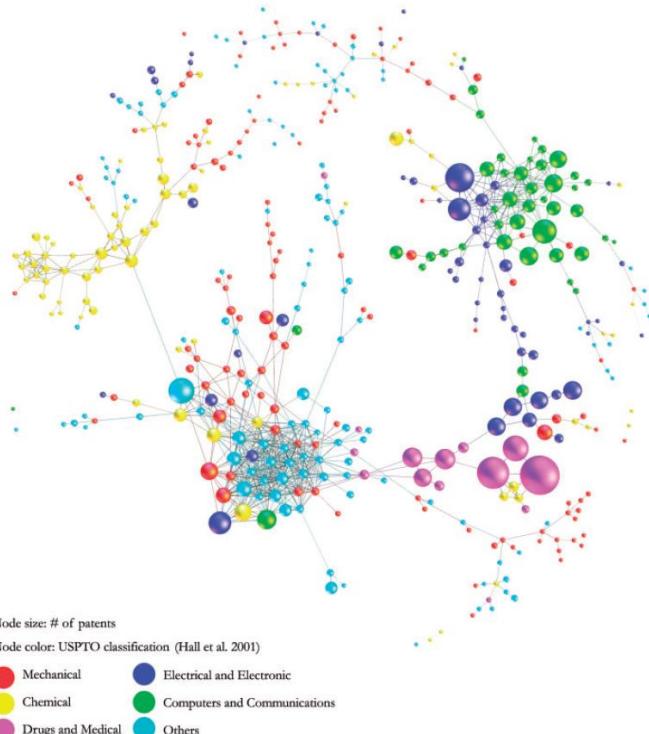
Evolutionary approaches in economic geography have contributed substantially to the growing body of knowledge of regional development processes and their underlying mechanisms. One key concept in the literature on evolutionary economic geography is that of *related variety*. Herein, regional industry structure is represented through the level of related variety of technologies, skills, or outputs. The related variety concept proposes that regional economic development is favored when an economy diversifies into products or technologies that are closely related to the stock of existing activities. In this article, we raise substantive questions regarding the internal logic of the concept of related variety, its spatial expressions, measurement specifics, empirical regularities and biases, and its possible short- and long-term effects on regional development. Based on this investigation, we make suggestions for improvements to future research.

ECONOMIC GEOGRAPHY

441

99(5)441–470 © 2023 Clark University.

[www.economicgeography.org](http://www.economicgeography.org)



"The most influential business thinker on Earth."  
—The New Yorker

# THE INNOVATOR'S DILEMMA

WHEN NEW TECHNOLOGIES CAUSE GREAT FIRMS TO FAIL //

CLAYTON M.  
CHRISTENSEN  
BESTSELLING AUTHOR OF  
THE INNOVATOR'S SOLUTION

HARVARD BUSINESS REVIEW PRESS

# Transformers and BDS

A BDS approach identifying novel combinations in firms (remember?)

## REPORTS

ring of likely core temperatures. However, further degenerations into multi-component systems are needed to fully understand their effect on the elastic properties of the core. Overall, our results demonstrate that the inner core has to be in the strongly nonisotropic state to have enough energy to invoke special circumstances such as strong anelasticity, partial melting, or combinations of crystalline phases in order to match the observed seismic velocities and densities of the inner core.

### References and Notes

1. J.-P. Bouchet, *Geophys. Res. Lett.* **43**, 4377–4380 (1996).
2. J.-P. Bouchet, *Phys. Earth Planet. Inter.* **85**, 119–137 (1994).
3. A. Cao, R. Homanen, K. Takemoto, *Science* **100**, 1443–1445 (2000).
4. J.-P. Bouchet, *Phys. Earth Planet. Inter.* **254**, 227–232 (2007).
5. K. Xie, K. E. Nelson, *Phys. Rev. E* **81**, 014502 (2005).
6. K. Xie, R. E. Cohen, *Geophys. Res. Lett.* **37**, 023025 (2010).
7. D. Annaswamy et al., *Earth Planet. Sci. Lett.* **225**, 243–253 (2004).
8. D. Annaswamy et al., *Phys. Earth Planet. Inter.* **164**, 63–89 (2007).
9. D. Annaswamy, J. Brodholt, J. G. Woodward, *Earth Planet. Sci. Lett.* **245**, 143–151 (2006).
10. M. A. Schwartz, D. L. Anderson, *Phys. Earth Planet. Inter.* **165**, 373–378 (2001).
11. D. Annaswamy et al., *Earth Planet. Sci. Lett.* **295**, 303–310 (2010).
12. Z. Wu et al., *Proc. Natl. Acad. Sci. U.S.A.* **109**, 10219–10224 (2012).
13. J.-P. Bouchet, *J. Appl. Phys.* **93**, 2472–2480 (2000).
14. J.-P. Bouchet, C. Rabaté, C. Ravel-Chapelle, *J. Appl. Phys.* **84**, 4445–4455 (2000).
15. J.-P. Bouchet, *Phys. Rev. E* **62**, 046110 (2000).
16. J.-P. Bouchet, J. Bourgat, D. Villedieu, S. Marrot, F. Garet, *Geophys. J. Int.* **145**, 501–509 (2001).
17. J.-P. Bouchet, *High-Density Phys. II*, 141–144 (2011).
18. D. Annaswamy, *Phys. Earth Planet. Inter.* **200**, 156–167 (2010).
19. Y. W. Galvin, D. J. Swanson, *J. Phys. Chem. Solids* **35**, 126–134 (1974).
20. D. R. Nelson, *Dynamical Geometry in Condensed Matter Physics* (Princeton Univ. Press, Princeton, NJ, 2002).
21. Y. W. Galvin, *J. Phys. B* **2**, 3952–3959 (1970).
22. D. R. Nelson, *Phys. Rev. E* **50**, 5253–5260 (1994).
23. D. R. Nelson, *Geophys. Res. Lett.* **21**, 1743–1746 (1994).
24. V. Saito, T. Pohjola, J. Alton, *Phys. Rev. E* **84**, 016105 (2011).
25. V. Saito, T. Pohjola, J. Alton, *Phys. Rev. E* **84**, 071130 (2011).
26. F. Delgaty, *Phys. Chem. Chem. Phys.* **8**, 3283–3292 (2006).
27. F. Delgaty, *J. Phys. Chem. B* **110**, 3283–3292 (2006).
28. F. Delgaty, *J. Phys. Chem. B* **110**, 3293–3293 (2006).
29. F. Delgaty, *J. Phys. Chem. B* **110**, 3293–3293 (2006).
30. G. M. Ward, F. Delgaty, *J. Mater. Sci.* **42**, 4673–4679 (2007).
31. G. M. Ward, F. Delgaty, *J. Mater. Sci.* **42**, 4673–4679 (2007).
32. G. M. Ward, *Introduction to the Physics of the Earth's Interior* (Cambridge Univ. Press, Cambridge, 2000), pp. 209–244.

**Acknowledgments.** Supported by National Environment Research Council (NERC) grants NE/C000520/1 and NE/C000521/1. The authors would like to thank the referees whose comments were helpful in the NERC-supervised chapter. Computer simulations were carried out on the NERSC supercomputer. This work was funded by the NERC and the University of Bristol. We thank the editor and the anonymous reviewers for their useful comments. L.W. designed research, analyzed data, and wrote the paper. J.B. analyzed data and wrote the paper. I.C.W. analyzed data and wrote the paper. D.R.N. analyzed data and wrote the paper.

### Supplementary Materials

[www.sciencemag.org/lookup/doi/10.1126/science.1246134](http://www.sciencemag.org/lookup/doi/10.1126/science.1246134)

Supplementary Text

Fig. S1 to S5

References

123–40

23 July 2013; accepted 26 September 2013

Published online 11 October 2013;

10.1126/science.1246134

between extending science with atypical combinations of knowledge and combining what the advantages of conventional domain-level thinking is critical to the link between innovativeness and impact. However, little is known about the composition of this balance or how scientists can achieve it.

In this study, we examined 179 million references in the bibliography of each paper (23, 24). We then calculated the frequency of each co-citation between papers in the WOS and compared these observed frequencies to those expected by chance, using randomized citation networks. In the randomized citation network, the probability of a citation between two papers in the WOS was selected by means of a Monte Carlo algorithm. The switching algorithm performed 1000 iterations to ensure that each paper and the distribution of these citation counts forward and backward in time to ensure that each paper (or journal) with  $n$  citations in the observed network had  $n$  citations in the randomized network. For both the observed and the randomized papers-paper citation networks, the respective journal pairs to focus on domain-level combinations (24–26). In the dataset, there were over 122 million unique journal pairs represented by the 15613 journals indexed in the WOS.

We considered pairwise combinations of references. We then calculated the frequency of each co-citation between papers in the WOS and compared these observed frequencies to those expected by chance, using randomized citation networks. For both the observed and the randomized papers-paper citation networks, the respective journal pairs to focus on domain-level combinations (24–26). In the dataset, there were over 122 million unique journal pairs represented by the 15613 journals indexed in the WOS.

## Atypical Combinations and Scientific Impact

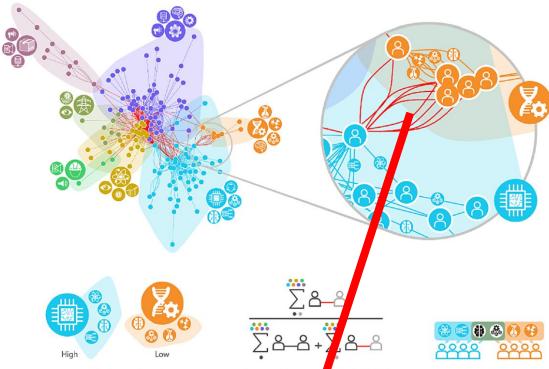
Brian Uzzi,<sup>1,2</sup> Satyam Mukherjee,<sup>1,2</sup> Michael Stringer,<sup>2,3</sup> Ben Jones,<sup>4,5\*</sup>

Novelty is an essential feature of creative ideas, yet the building blocks of new ideas are often embodied in existing knowledge. From this perspective, balancing atypical knowledge with conventional knowledge may be critical to the link between innovativeness and impact. Our analysis of 179 million scientific publications in the Web of Science (WOS) found that current interest in team science and scientific impact is premised on the idea that the best ideas are born from the combination of multiple experts with different knowledge. Yet combining knowledge that prompts scientific breakthroughs (*i.e.*, “transformers”)

Scientific enterprises are increasingly constrained to research within narrow boundaries. It is unlikely to be the source of the most fruitful ideas (*i.e.*, Models of creativity emphasize that individual scientists inevitably become narrower in their expertise as the breadth of their knowledge increases). As a result, creative breakthroughs across borders will be increasingly challenging (*i.e.*, especially given the difficulty of identifying the right partners). Yet, as a general rule, novel ideas can be difficult to absorb (*i.e.*, and communicate, leading scientists to initially dismiss them). In *The Pragmatic Guide to Creativity*, Christensen et al. argue that the lack of appreciation for the need to absorb and communicate knowledge from outside one’s field, and the resulting reluctance to do so, is a major impediment to breakthroughs (*i.e.*, “BDS”)

<sup>1</sup> Kellogg School of Management, Northwestern University, 2001 Sheridan Road, Evanston, IL 60202, USA; <sup>2</sup> Northeastern Institute of Complex Systems, Northeastern University, 360 Brattle Street, Cambridge, MA 02134, USA; <sup>3</sup> Department of Biological Sciences, DePaul University, 2550 North Adams Street, Chicago, IL 60614, USA; <sup>4</sup> National Bureau of Economic Research, 1050 Massachusetts Avenue, Cambridge, MA 02138, USA

\*Corresponding author. E-mail: bennettjones@kellogg.northwestern.edu

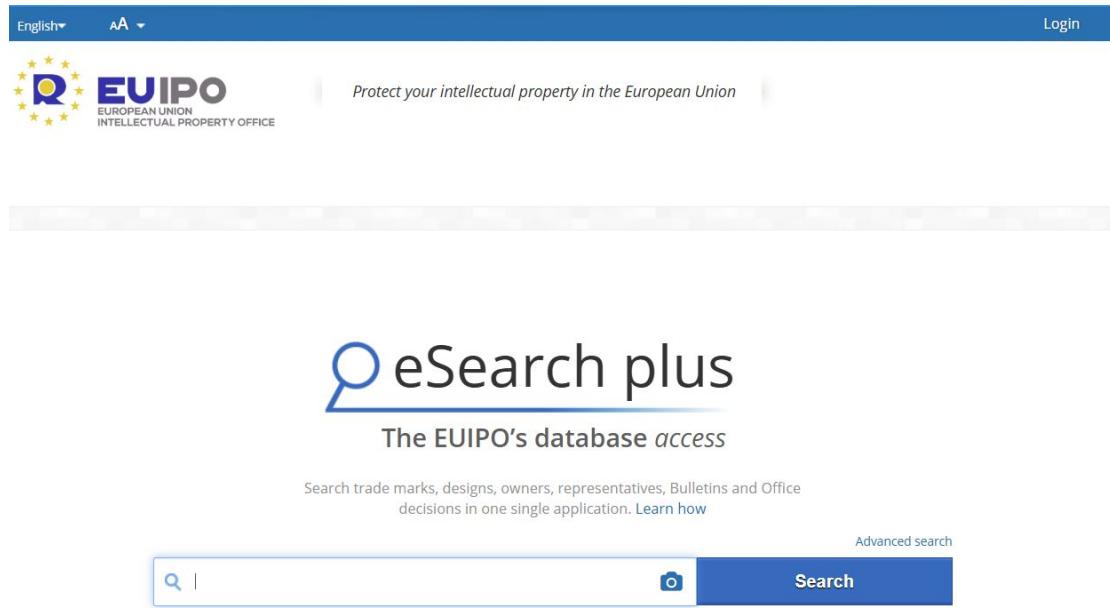


## **Market2Vec (or MarketBert)**

Can we predict a firm's next products (or product features) based on its past products?

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?



The screenshot shows the homepage of the EUIPO eSearch plus database. At the top, there is a blue header bar with 'English' and 'AA' language options, and a 'Login' button. Below the header is the EUIPO logo, which includes a stylized 'Q' made of yellow stars and the text 'EUIPO' in blue, with 'EUROPEAN UNION INTELLECTUAL PROPERTY OFFICE' underneath. To the right of the logo is the tagline 'Protect your intellectual property in the European Union'. The main title 'eSearch plus' is displayed prominently in large, dark blue font, with a magnifying glass icon integrated into the letter 'e'. Below the title, the subtitle 'The EUIPO's database access' is shown. A sub-subtitle explains the functionality: 'Search trade marks, designs, owners, representatives, Bulletins and Office decisions in one single application. Learn how'. There is also a link to 'Advanced search'. At the bottom, there is a search bar with a magnifying glass icon and a camera icon, followed by a blue 'Search' button.

<https://euipo.europa.eu/eSearch/>

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?

<input type="checkbox"/> 001765916 - novo nordisk <a href="#">+info</a>	<input type="checkbox"/> 015703441 - novo nordisk <a href="#">+info</a>	<input type="checkbox"/> 018213481 - novo nordisk <a href="#">+info</a>
 <b>Trade mark information</b> Trade mark number <b>001765916</b> Type <b>Figurative</b> Filing date <b>20/07/2000</b> Registration date <b>20/11/2001</b> Nice Classification <b>5, 9, 10, 16, 41, 42</b> Trade mark status <b>Registered</b> Basis <b>EUTM</b> Reference <b>200285</b>	 <b>Trade mark information</b> Trade mark number <b>015703441</b> Type <b>Figurative</b> Filing date <b>28/07/2016</b> Registration date <b>24/11/2016</b> Nice Classification <b>1, 5</b> Trade mark status <b>Registered</b> Basis <b>EUTM</b> Reference <b>11971/EM/JVe</b>	 <b>Trade mark information</b> Trade mark number <b>018213481</b> Type <b>Figurative</b> Filing date <b>20/03/2020</b> Registration date <b>08/09/2020</b> Nice Classification <b>1, 5, 9, 10, 41, 42, 44</b> Trade mark status <b>Registered</b> Basis <b>EUTM</b> Reference <b>12918/EM</b>

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?

001765916 - novo nordisk [+info](#)



Trade mark information

Trade mark number	001765916
Type	Figurative
Filing date	20/07/2000
Registration date	20/07/2001
Nice Classification	5, 9, 10, 16, 41, 42
Trade mark status	Registered
Basis	EUTM
Reference	200285

018213481 - novo nordisk [+info](#)



Trade mark information

Trade mark number	018213481
Type	Figurative
Filing date	20/03/2020
Registration date	08/09/2020
Nice Classification	1, 5, 9, 10, 41, 42, 44
Trade mark status	Registered
Basis	EUTM
Reference	12918/EM

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?

Filing date	20/07/2000		Filing date	20/03/2020	
Registration date	20/11/2001		Registration date	08/09/2020	
Nice Classification	5, 9, 10, 16, 41, 42		Nice Classification	1, 5, 9, 10, 41, 42, 44	

**5** Pharmaceutical preparations and substances.

**9** Recorded computer software, including software for medical information, education and sales activities, software for websites, intranet and Internet; medical publications in electronic form; data processing equipment and computers.

**10** Medical and surgical apparatus and instruments; syringes and needles for medical purposes; parts and fittings (not included in other classes) for all the aforementioned goods.

**16** Periodical publications; printed matter, instructional and teaching material (except apparatus) all for medical purposes, software manuals.

**41** Medical training and education, IT training and education.

**42** Computer services, namely providing access to data concerning medical analysis and treatment; medical counselling services, medical clinical services, medical services related to diabetes preparations and other pharmaceuticals; information and advisory services all relating to personal management of medical conditions; providing IT services and IT consultancy, including development, implementation and maintenance of IT system solutions.

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?



**5** Pharmaceutical preparations and substances.

**9** Recorded computer software, including software for medical information, education and sales activities, software for websites, intranet and Internet; medical publications in electronic form; data processing equipment and computers.

**10** Medical and surgical apparatus and instruments; syringes and needles for medical purposes; parts and fittings (not included in other classes) for all the aforementioned goods.

**16** Periodical publications; printed matter, instructional and teaching material (except apparatus) all for medical purposes, software manuals.

**41** Medical training and education, IT training and education.

**42** Computer services, namely providing access to data concerning medical analysis and treatment; medical counselling services, medical clinical services, medical services related to diabetes preparations and other pharmaceuticals; information and advisory services all relating to personal management of medical conditions; providing IT services and IT consultancy, including development, implementation and maintenance of IT system solutions.

**1** Chemical preparations for use in industry; chemical additives; chemical additives for use in the manufacture of pharmaceuticals, cosmetics, hair lotions, disinfectants; chemical diagnostic reagents for industrial use; chemical preparations for use in the oil industry.

**44** Providing medical services via websites and mobile apps for the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment; providing medical information and medical counselling the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment.

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?



**5** Pharmaceutical preparations and substances.

**9** Recorded computer software, including software for medical information, education and sales activities, software for websites, intranet and Internet; medical publications in electronic form; data processing equipment and computers.

**10** Medical and surgical apparatus and instruments; syringes and needles for medical purposes; parts and fittings (not included in other classes) for all the aforementioned goods.

**16** Periodical publications; printed matter, instructional and teaching material (except apparatus) all for medical purposes, software manuals.

**41** Medical training and education, IT training and education.

**42** Computer services, namely providing access to data concerning medical analysis and treatment; medical counselling services, medical clinical services, medical services related to diabetes preparations and other pharmaceuticals; information and advisory services all relating to personal management of medical conditions; providing IT services and IT consultancy, including development, implementation and maintenance of IT system solutions.

**1** Chemical preparations for use in industry; chemical additives; chemical additives for use in the manufacture of pharmaceuticals, cosmetics, hair lotions, disinfectants, chemical diagnostic reagents for industrial use; chemical preparations for use in the oil industry.

**44** Providing medical services via websites and mobile apps for the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment; providing medical information and medical counselling the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment.

Concept (or market) token?

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?



5 Pharmaceutical preparations and substances.

9 Recorded computer software, including software for medical information, education and sales activities, software for websites, intranet and Internet; medical publications in electronic form; data processing equipment and computers.

10 Medical and surgical apparatus and instruments; syringes and needles for medical purposes; parts and fittings (not included in other classes) for all the aforementioned goods.

16 Periodical publications; printed matter, instructional and teaching material (except apparatus) all for medical purposes, software manuals.

41 Medical training and education, IT training and education.

42 Computer services, namely providing access to data concerning medical analysis and treatment; medical counselling services, medical clinical services, medical services related to diabetes preparations and other pharmaceuticals; information and advisory services all relating to personal management of medical conditions; providing IT services and IT consultancy, including development, implementation and maintenance of IT system solutions.

1 Chemical preparations for use in industry; chemical additives; chemical additives for use in the manufacture of pharmaceuticals, cosmetics, hair lotions, disinfectants, chemical diagnostic reagents for industrial use; chemical preparations for use in the oil industry.

44 Providing medical services via websites and mobile apps for the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment; providing medical information and medical counselling the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment.

Concept (or market) token?

Challenge!!  
How to construct market token sequences?

# Market2Vec (or MarketBert)

Can we predict a firm's next products (or product features) based on its past products?



5 Pharmaceutical preparations and substances.

9 Recorded computer software, including software for medical information, education and sales activities, software for websites, intranet and Internet; medical publications in electronic form; data processing equipment and computers.

10 Medical and surgical apparatus and instruments; syringes and needles for medical purposes; parts and fittings (not included in other classes) for all the aforementioned goods.

16 Periodical publications; printed matter, instructional and teaching material (except apparatus) all for medical purposes, software manuals.

41 Medical training and education, IT training and education.

42 Computer services, namely providing access to data concerning medical analysis and treatment; medical counselling services, medical clinical services, medical services related to diabetes preparations and other pharmaceuticals; information and advisory services all relating to personal management of medical conditions; providing IT services and IT consultancy, including development, implementation and maintenance of IT system solutions.

1 Chemical preparations for use in industry; chemical additives; chemical additives for use in the manufacture of pharmaceuticals, cosmetics, hair lotions, disinfectants, chemical diagnostic reagents for industrial use; chemical preparations for use in the oil industry.

44 Providing medical services via websites and mobile apps for the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment; providing medical information and medical counselling the management of diabetes, obesity, haemostasis management and for the treatment of sickle cell disease, treatment of brain disorders, non-alcoholic steatohepatitis, atherosclerosis, cardiovascular disease, chronic kidney disease, hormone therapy and stem cell treatment.

Concept (or market) token?

Challenge!!  
How to construct market token sequences?

# Market2Vec (o

Can we predict a firm's next products  
past products?

Filing date	20/07/2000
Registration date	20/07/2001
Nice Classification	5, 9, 10, 16, 41, 42

- 5 Pharmaceutical preparations and substances.
- 9 Recorded computer software, including software for medical information, education and sales activities, software for websites, Internet; medical publications in electronic form; data processing equipment and computers.
- 10 Medical and surgical apparatus and instruments; syringes and needles for medical purposes; parts and fittings (not included in classes) for all the aforementioned goods.
- 16 Periodical publications; printed matter, instructional and teaching material (except apparatus) all for medical purposes, software.
- 41 Medical training and education, IT training and education.
- 42 Computer services, namely providing access to data concerning medical analysis and treatment; medical counselling services, medical clinical services, medical services related to diabetes preparations and other pharmaceuticals; information and advisory services all relating to personal management of medical conditions; providing IT services and IT consultancy, including development, implementation and maintenance of IT system solutions.

1 Chemical preparations for use in industry; chemical hair lotions, disinfectants, chemical diagnostic reagents

44 Providing medical services via web treatment of sickle cell disease, treatment of kidney disease, hormone therapy and treatment of obesity, haemostasis management and treatment of atherosclerosis, cardiovascular disease, i

Concept (or market) token?

```
# =====
# 4) TRAINING LOOP (4 steps: forward -> loss -> backward -> update)
# =====
def train_one_epoch(model, train_loader, optimizer, scheduler):
    model.train()
    total = 0.0
    pbar = tqdm(train_loader, desc="Train", leave=False)

    for batch in pbar:
        batch = {k: v.to(device) for k, v in batch.items()}

        # (1) Forward pass
        outputs = model(**batch)

        # (2) Loss
        loss = outputs.loss

        # (3) Backward
        loss.backward()
        torch.nn.utils.clip_grad_norm_(model.parameters(), GRAD_CLIP)

        # (4) Update
        optimizer.step()
        scheduler.step()
        optimizer.zero_grad()

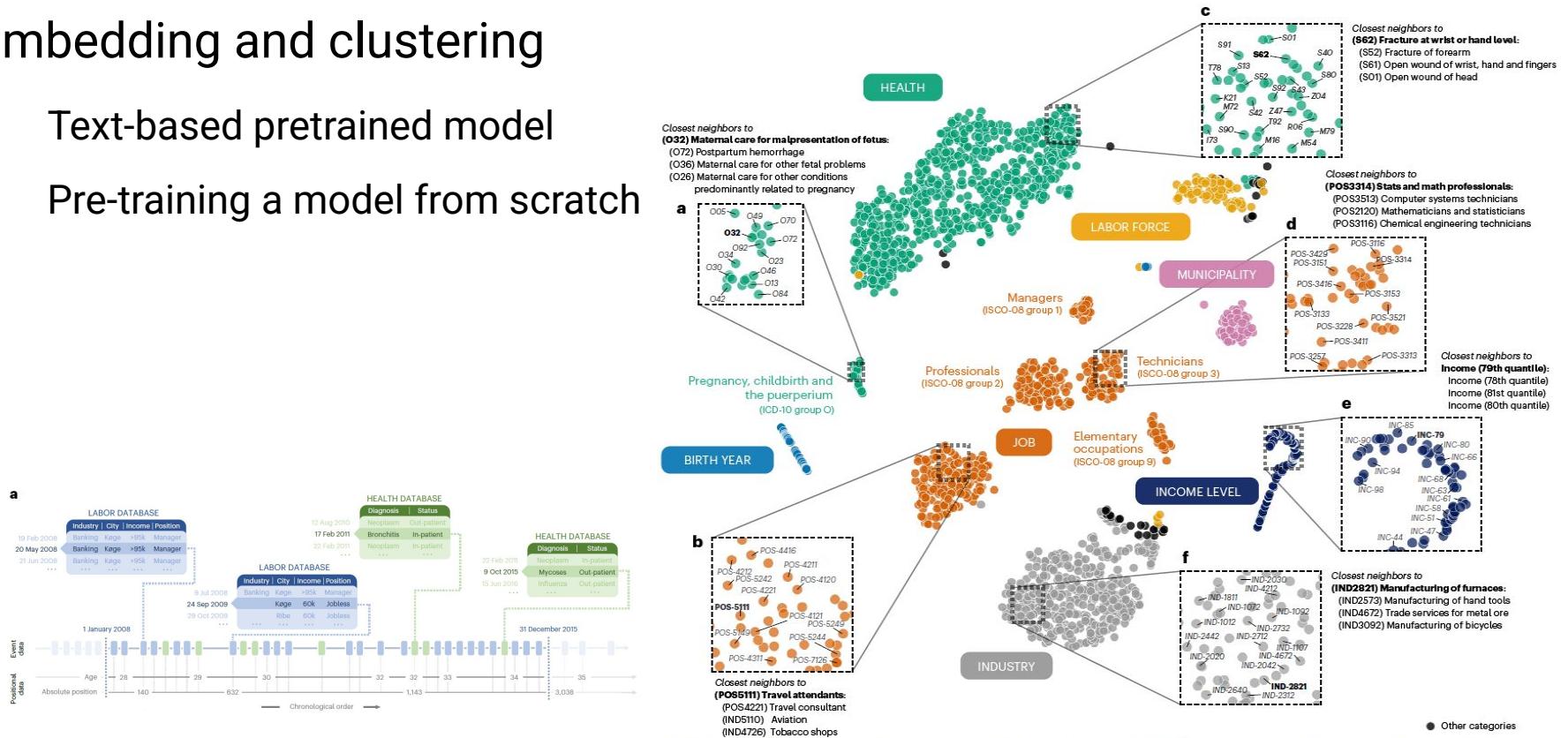
        total += loss.item()
        pbar.set_postfix(loss=f"{loss.item():.4f}")

    return total / max(1, len(train_loader))
```

# Market2Vec (or MarketBert)

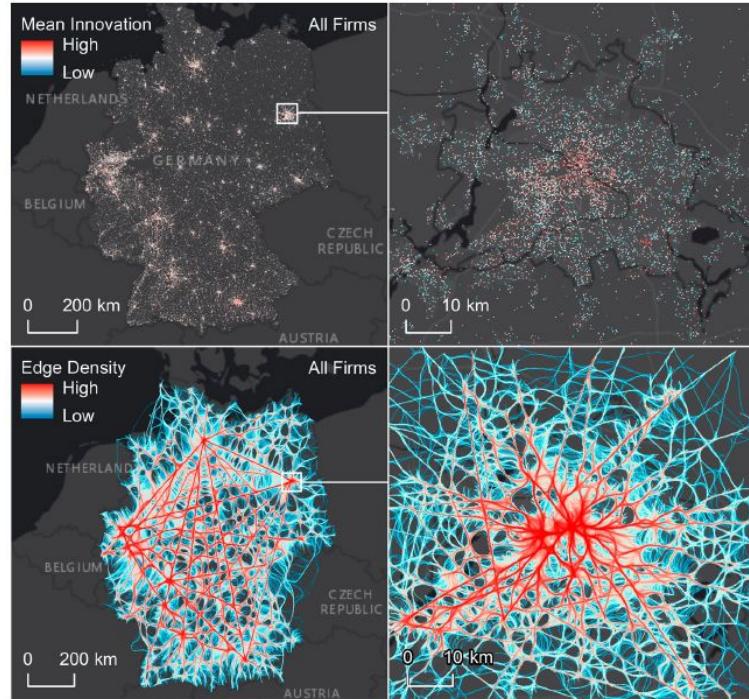
# Embedding and clustering

- Text-based pretrained model
  - Pre-training a model from scratch



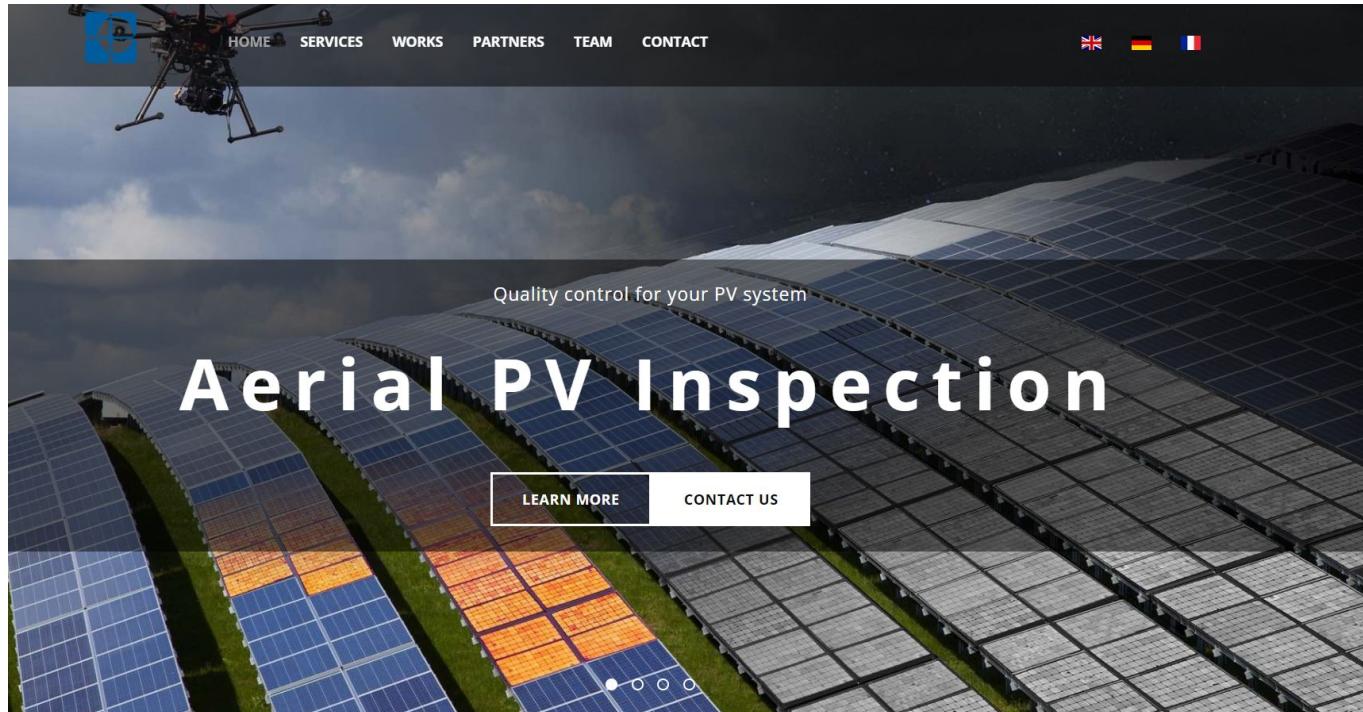
# Sentence Transformers Finetuning (SetFit)

Text-based pretrained model (using SetFit for finetuning)



# Sentence Transformers Finetuning (SetFit)

Text-based pretrained model (using SetFit for finetuning)



# Sentence Transformers Finetuning (SetFit)

Text-based pretrained model (using SetFit for finetuning)

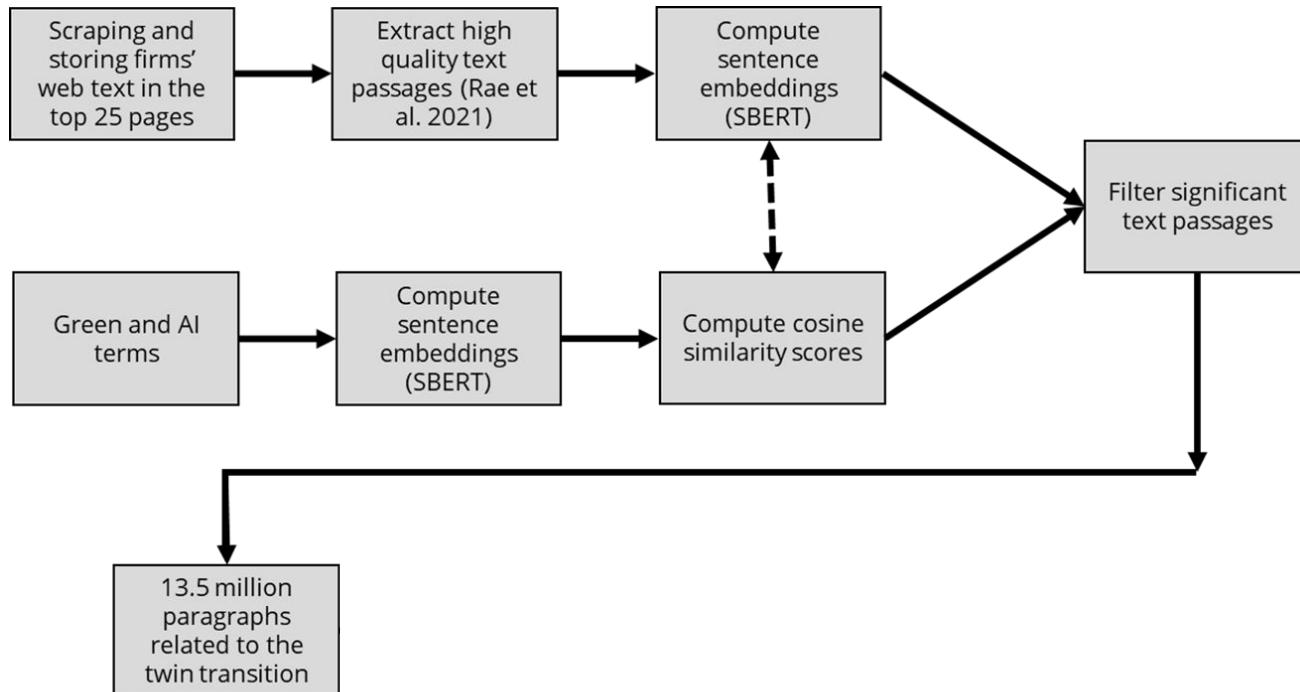


We examine the quality of photovoltaic systems without dismantling.  
Are the modules ok? Were they damaged during assembly? Are the yields too low?  
Aerial PV inspection identifies and locates faults precisely.

We combine and automate classic and preventive inspections of solar power plants with image-based methods like thermography, electroluminescence and UV-fluorescence measurements, and with IV-curve recording by day and night, isolation issues detection, serial number scan, and geography data referencing.

To this end we employ tripod- and drone-based sensor systems, providing a cost-effective check-up by systematic sampling or 100% screening. EL measurements offer a comprehensive and exact assessment of faults on cell level and thus represent the tower of strength of the AePVI concept.

# Sentence Transformers Finetuning (SetFit)



# Sentence Transformers Finetuning (SetFit)

## Efficient Few-Shot Learning Without Prompts

Lewis Tunstall<sup>1</sup>, Nils Reimers<sup>2</sup>, Unso Eun Seo Jo<sup>1</sup>, Luke Bates<sup>3</sup>, Daniel Korat<sup>4</sup>, Moshe Wasserblat<sup>4</sup>, Oren Pereg<sup>\*</sup>

<sup>1</sup>Hugging Face <sup>2</sup>cohere.ai

<sup>3</sup>Ubiquitous Knowledge Processing Lab, Technical University of Darmstadt

<sup>4</sup>Emergent AI Lab, Intel Labs

<sup>1</sup>firstname@huggingface.com <sup>2</sup>info@nils-reimers.de

<sup>3</sup>bates@ukp.informatik.tu-darmstadt.de

<sup>4</sup>firstname.lastname@intel.com

## Abstract

Recent few-shot methods, such as parameter-efficient fine-tuning (PEFT) and pattern-exploring training (PET), have achieved impressive results in label-scarce settings. However, they are difficult to employ since they are subject to high variability from manually crafted prompts and typically require billion-parameter language models to achieve high accuracy. To address these shortcomings, we propose SetFit (Sentence Transformer Finetuning), an efficient and prompt-free framework for few-shot fine-tuning of Sentence Transformers (ST). SetFit works by first finetuning a pretrained ST on a small number of text pairs, in a contrastive Siamese manner. The resulting model is then used to generate rich text embeddings, which are used to train a classification head. This simple framework requires no prompts or verbalizers, and achieves high accuracy with orders of magnitude less parameters than existing techniques. Our experiments show that SetFit obtains comparable results with PEFT and PET techniques, while being an order of magnitude faster to train. We also show that SetFit can be applied in multilingual settings by simply switching the ST body. Our code<sup>1</sup> and datasets<sup>2</sup> are made publicly available.

## 1 Introduction

Few-shot learning methods have emerged as an attractive solution to label-scarce scenarios, where data annotation can be time-consuming and costly. These methods are designed to work with a small number of labeled training examples, and typically involve adapting pretrained language models (PLMs) for specific downstream tasks.

Today, there exist several approaches to few-shot learning with PLMs. These include in-context learning (ICL), parameter-efficient finetuning (PEFT), and pattern-exploiting training

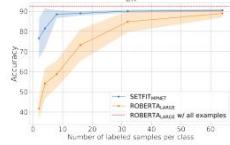
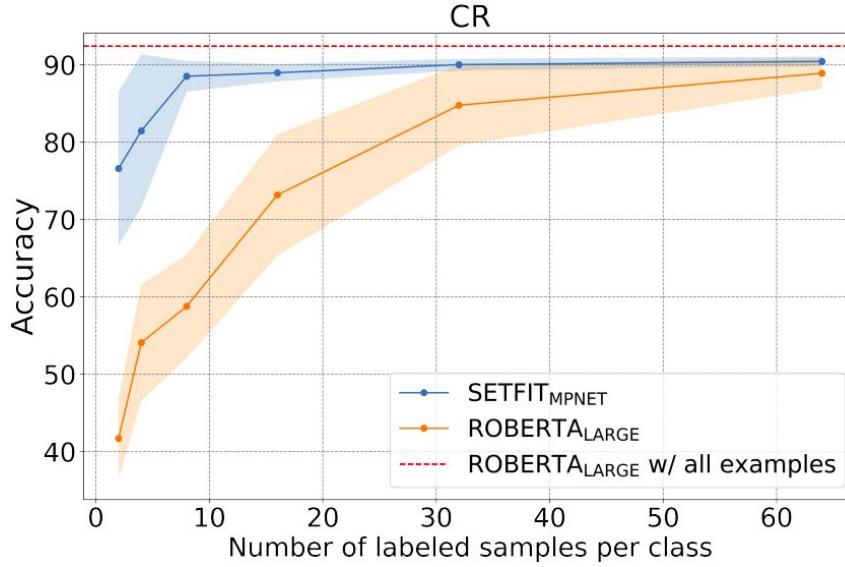


Figure 1: Compared to standard fine-tuning, SetFit is more sample efficient and exhibits less variability when trained on a small number of labeled examples.

(PET). Unfortunately, these approaches can be impractical for many researchers and practitioners. One disadvantage is that these approaches typically rely on the use of large-scale language models to achieve high performance. For example, T-FEW (Liu et al., 2022) is based on the 11 billion parameter model T0 (Sanh et al., 2021), while GPT-3 (Brown et al., 2020a) is an order of magnitude larger. Secondly, training and deploying these few-shot methods typically requires specialized infrastructure with limited accessibility. Moreover, PET and the prominent PEFT methods require, as part of their training, the input of manually generated prompts, yielding varying outcomes depending on the level of manual prompt-engineering.

In this paper, we propose SetFit, an approach based on Sentence Transformers (ST) (Reimers and Gurevych, 2019) that dispenses with prompts altogether and does not require large-scale PLMs to achieve high accuracy. For example, with only 8 labeled examples in the Customer Reviews (CR) sentiment dataset, SetFit is competitive with finetuning on the full training set, despite the fine-tuned model being three times larger (see Figure 1).

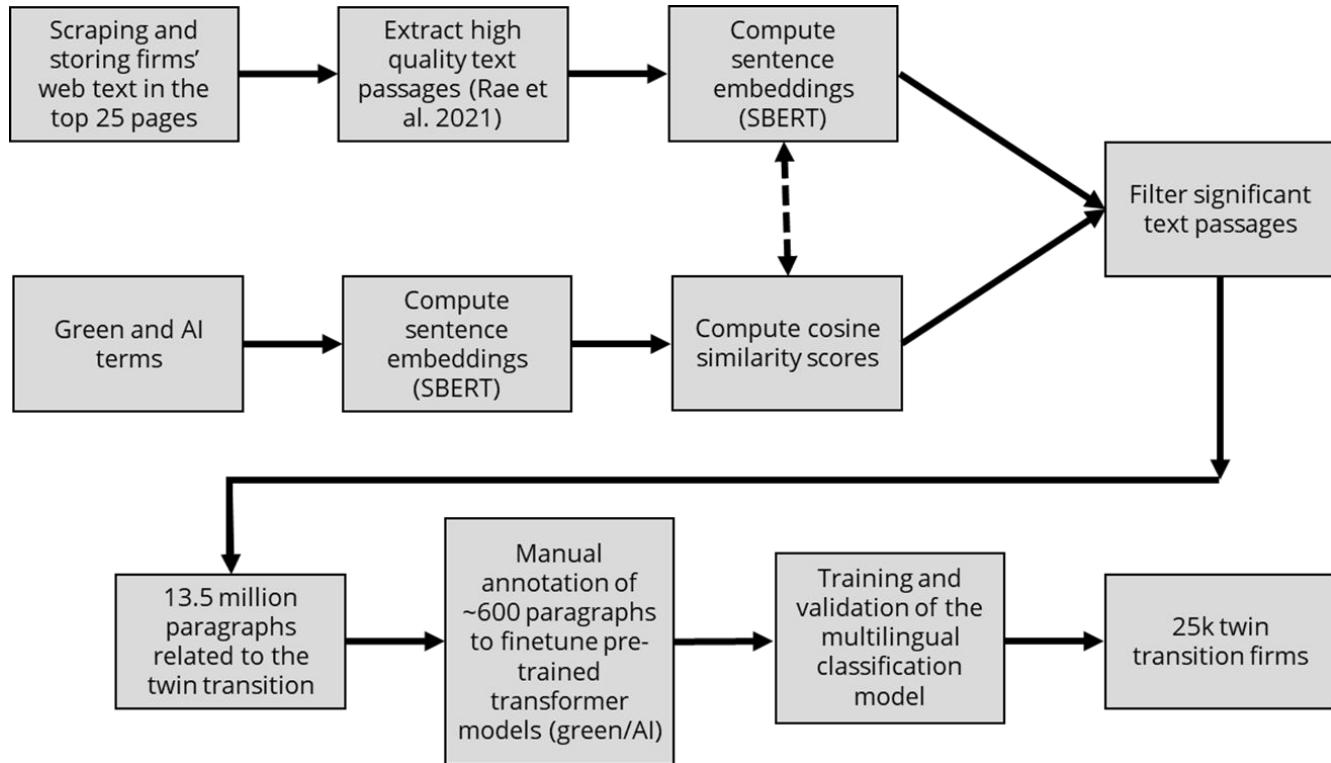


We demonstrate SetFit's efficacy in few-shot text classification over a range of NLP datasets

<sup>1</sup><https://github.com/huggingface/setfit>

<sup>2</sup><https://huggingface.co/setfit>

# Market2Vec (or MarketBert)



# Sentence Transformers Finetuning (SetFit)

## Text-based pretrained model (using SetFit for finetuning)

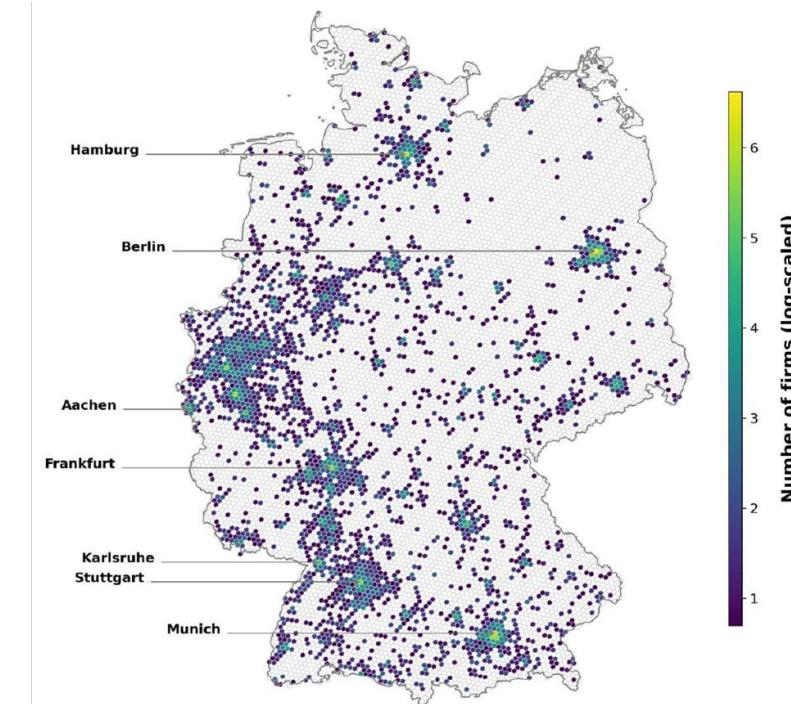
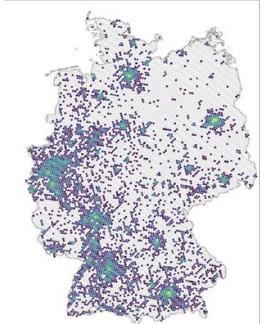
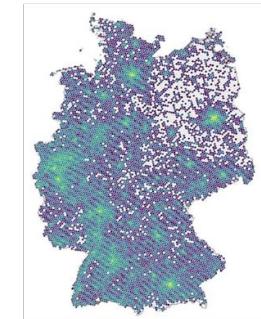


The twin transition, which refers to the combined shift of economies toward greener and more digital modes of production, is highly prioritised on the political agenda in Europe

**CONTACT** Sebastian Losacker [sebastian.losacker@geogru-gießen.de](mailto:sebastian.losacker@geogru-gießen.de)  
Geography and Geoinformatics, Institute of Spatial Sciences, University of Göttingen, Germany  
Department of Economic Geography, Faculty of Social Sciences, University of Groningen, the Netherlands  
CIRCLE – Centre for Innovation Research, Lund University, Sweden

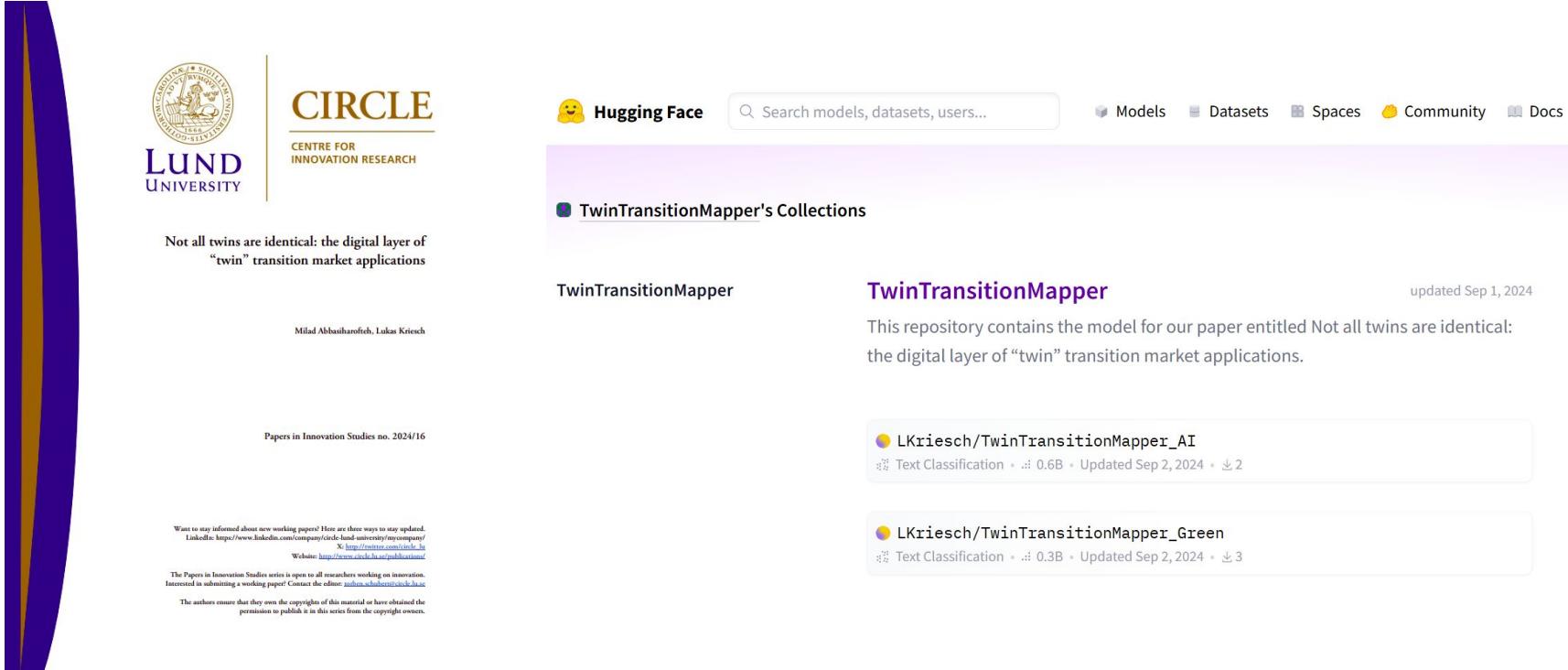
Supplemental data for this article can be accessed online at <https://doi.org/10.1080/21681376.2025.2510679>

© 2025 The Authors. Published by Informa UK Limited, trading as Taylor & Francis Group.  
This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. The specific terms on which this article has been published allow the posting of the Accepted Manuscript in a repository by the author(s) or their institution.



# Sentence Transformers Finetuning (SetFit)

Text-based pretrained model (using SetFit for finetuning)





LUND  
UNIVERSITY



CIRCLE  
CENTRE FOR  
INNOVATION  
RESEARCH

Not all twins are identical: the digital layer of  
“twin” transition market applications

Milad Abbasifarreh, Lukas Kriesch

Papers in Innovation Studies no. 2024/16

Want to stay informed about new working papers? Here are three ways to stay updated.  
LinkedIn: <https://www.linkedin.com/company/circle-lund-university-research/>  
X: [https://twitter.com/circle\\_lu](https://twitter.com/circle_lu)  
Website: <http://www.circle.lu/en/publications/>

The Papers in Innovation Studies series is open to all researchers working on innovation  
Interested in submitting a working paper? Contact the editor: [authors@circle.lu/en/](mailto:authors@circle.lu/en/)

The authors ensure that they own the copyrights of this material or have obtained the  
permission to publish it in this series from the copyright owners.



Hugging Face

Search models, datasets, users...

Models Datasets Spaces Community Docs



TwinTransitionMapper's Collections

TwinTransitionMapper

TwinTransitionMapper

updated Sep 1, 2024

This repository contains the model for our paper entitled Not all twins are identical:  
the digital layer of “twin” transition market applications.



LKriesch/TwinTransitionMapper\_AI

Text Classification · 0.6B · Updated Sep 2, 2024 · ↴ 2



LKriesch/TwinTransitionMapper\_Green

Text Classification · 0.3B · Updated Sep 2, 2024 · ↴ 3

# **Going beyond Market2Vec**

Going beyond Market2Vec as a prediction tool

- What does the bias in a transformer model tell us about firms' next potential products?
- Can we interpret lower levels of predictability in our model, in some cases, as meaningful signals rather than mere noise?

## Going beyond Market2Vec

Applying the same approach in a broader context:

- Emerging skills in the job market
- Detecting scientific breakthrough
- Atypical inventions ☐ firm's future market value!
- Investigating team compositions
- etc.