

**dfki** Deutsches Forschungszentrum für Künstliche Intelligenz German Research Center for Artificial



UNIVERSITÄT DES SAARLANDES

Intelligence

# **Large Process Models**

Peter Fettke

18th Symposium and Summer School on Service-Oriented Computing (SummerSoC) Crete, Greece June 24th - 29th, 2024



Deutsches Forschungszentrum für Künstliche Intelligenz German Research Center for Artificial



UNIVERSITÄT DES SAARLANDES

Intelligence

# **Large Process Models**

Peter Fettke

in cooperation with Wolfgang Reisig

HUMBOLDT-UNIVERSITÄT ZU BERLIN



18th Symposium and Summer School on Service-Oriented Computing (SummerSoC) Crete, Greece June 24th - 29th, 2024

#### This talk in a nutshell



#### This talk in a nutshell



#### Agenda

dfk

Deutsches Forschungszentrum für Künstliche Intelligenz German Research Center for Artificial



UNIVERSITÄT DES SAARLANDES

Intelligence

Large Language Models
 Process Models
 Large Process Models

HUMBOLDT-UNIVERSITÄT ZU BERLIN



## What is a large language model? Most abstract view











10

## What is a large language model? A technical view



## What is a large language model? A technical view



What is a token?

- character
- word
- n-gram
- part of a word
- byte pair encoding (BPE)









#### Article Accurate structure prediction of biomolecular interactions with AlphaFold 3

os://doi.org/10.1038/s41586-024-07487-w	Josh Abramson <sup>1,7</sup> , Jonas Adler <sup>1,7</sup> , Jack Dunger <sup>1,7</sup> , Richard Evans <sup>1,7</sup> , Tim Green <sup>1,7</sup> ,	
ceived: 19 December 2023	Alexander Pritzel <sup>47</sup> , Olaf Ronneberger <sup>47</sup> , Lindsay Willinore <sup>32</sup> , Andrew J. Ballard <sup>4</sup> , Joshua Bambrick <sup>5</sup> , Sebastian W. Sofenstein <sup>7</sup> , David A. Evans <sup>2</sup> , Chi-Chun Hung <sup>3</sup> , Michael O'Neill <sup>5</sup> , David Beiman <sup>7</sup> , Kathyn Tumyasuvunakool <sup>7</sup> , Zachary Wu <sup>3</sup> , Akvillé Zamgutyté Erini Alvanili <sup>8</sup> , Charles Beatic <sup>5</sup> , Otavia Bentoll <sup>17</sup> , Alexa Bridgland <sup>7</sup> , Bara Bridgland <sup>7</sup> , Alexa Bridgland <sup>7</sup>	
cepted: 29 April 2024		
olished online: 8 May 2024		
en access		
Check for updates		
	Michal Zielinski", Augustin Zidek", Victor Bapst ", Pushmeet Konu", Max Jaderberg """, Demis Hassabis <sup>12,8</sup> & John M. Jumper <sup>18</sup>	

The introduction of AlphaFold 2<sup>1</sup> has spurred a revolution in modelling the structure of proteins and their interactions, enabling a huge range of applications in protein modelling and design<sup>2-6</sup>. Here we describe our AlphaFold 3 model with a substantially undated diffusion-based architecture that is capable of predicting the joint structure of complexes including proteins, nucleic acids, small molecules, ions and modified residues. The new AlphaFold model demonstrates substantially improved accuracy over many previous specialized tools: far greater accuracy for protein-ligand interactions compared with state-of-the-art docking tools, much higher accuracy for protein-nucleic acid interactions compared with nucleic-acid-specific predictors and substantially higher antibody-antigen prediction accuracy compared with AlphaFold-Multimer v.2.378. Together, these results show that high-accuracy modelling across biomolecular space is possible within a single unified deep-learning framework.

structure of protein-protein interactions.

Here we present AlphaFold 3 (AF3)-a model that is capable of

category, it achieves a substantially higher performance than strong

This is achieved by a substantial evolution of the AF2 architec-

ture and training procedure (Fig. 1d) both to accommodate more

general chemical structures and to improve the data efficiency of

learning. The system reduces the amount of multiple-sequence

alignment (MSA) processing by replacing the AF2 evoformer with

predicts the raw atom coordinates with a diffusion module, replace

ing the AF2 structure module that operated on amino-acid-specific

work to improve local structure) also enable us to eliminate

Accurate models of biological complexes are critical to our understanding of cellular functions and for the rational design of therapeutics<sup>2-4,9</sup>. Enormous progress has been achieved in protein structure types present in the Protein Data Bank<sup>32</sup> (PDB) (Fig. 1a,b). In all but one prediction with the development of AlphaFold<sup>1</sup>, and the field has grown tremendously with a number of later methods that build on methods that specialize in just the given task (Fig. 1c and Extended the ideas and techniques of AlphaFold 2 (AF2)<sup>10-12</sup>. Almost immediately Data Table 1), including higher accuracy at protein structure and the after AlphaFold became available, it was shown that simple input modifications would enable surprisingly accurate protein interaction predictions13-15 and that training AF2 specifically for protein inter action prediction vielded a highly accurate system7.

http Rec Ace

Put Op

These successes lead to the question of whether it is possible to accurately predict the structure of complexes containing a much wider range of biomolecules including ligands ions nucleic acids and modified residues, within a deep-learning framework. A wide range of predictors for various specific interaction types has been developed16-28, as well as one generalist method developed concurrently with the frames and side-chain torsion angles (Fig. 2b). The multiscale present work<sup>29</sup>, but the accuracy of such deep-learning attempts has nature of the diffusion process (low noise levels induce the netbeen mixed and often below that of physics-inspired methods<sup>30,31</sup>. Almost all of these methods are also highly specialized to particular stereochemical losses and most special handling of bonding patinteraction types and cannot predict the structure of general biomo terns in the network, easily accommodating arbitrary chemical lecular complexes containing many types of entities.

Core Contributor, Google DeepMind, London, UK, <sup>2</sup>Core Contributor, Isomorphic Labs, London, UK, <sup>3</sup>Google DeepMind, London, UK, <sup>4</sup>Isomorphic Labs, London, UK, <sup>4</sup>Department of Molecula Construction, Stanford University, Stanford University, Stanford, CA, USA, "Department of Computer Science, Princeton Null USA, "Department Science, Princeton Null USA, "Department Science, Princeton Null VI, USA, " equally: Josh Abrams Max Jaderberg, Demis Hassabis, John M. Jumper. We-mail: jaderberg@isomorphiclabs.com; dhcontact@google.com; jumper@google.com

components.

Nature | Vol 630 | 13 June 2024 | 493





used tokens:

- 20 common amino acid types
- an unknown type
- a gap token
- a mask token

Source:

Jumper et al.: Highly accurate protein structure prediction with AlphaFold. In: Nature, 2021, Vol. 596, pp. 583-589

## What is a large language model? A technical view (refined)



What is a token?

- character
- word
- n-gram
- part of a word
- byte pair encoding (BPE)
- amino acid types
- program code
- visual chunks

## What is a large language model? A technical view (refined)



What is a token?

- character
- word
- n-gram
- part of a word
- byte pair encoding (BPE)
- amino acid types

#### What is a *token* in case of *processes*?

#### Agenda

Large Language Models
 Process Models
 Large Process Models

Deutsches Forschungszentrum für Künstliche Intelligenz German Research Center for Artificial Intelligence



ð,

d

UNIVERSITÄT DES SAARLANDES

HUMBOLDT-UNIVERSITÄT ZU BERLIN





















four steps describe behavior



#### A little bit larger



bakery run =<sub>def</sub> bake • supply to aide • move to shop • sell

bakery run =<sub>def</sub> bake • supply to aide • move to shop • sell

(1)-bakery runs =<sub>def</sub> bakery run

bakery run =<sub>def</sub> bake • supply to aide • move to shop • sell

- (1)-bakery runs =<sub>def</sub> bakery run
- (2)-bakery runs  $=_{def}$  (1)-bakery runs bakery run

bakery run =<sub>def</sub> bake • supply to aide • move to shop • sell

- (1)-bakery runs =<sub>def</sub> bakery run
  (2)-bakery runs =<sub>def</sub> (1)-bakery runs bakery run
- (3)-bakery runs  $=_{def}$  (2)-bakery runs bakery run

bakery run =<sub>def</sub> bake • supply to aide • move to shop • sell

(1)-bakery runs	= <sub>def</sub>	bakery run
(2)-bakery runs	= <sub>def</sub>	(1)-bakery runs • bakery run
(3)-bakery runs	= <sub>def</sub>	(2)-bakery runs • bakery run

(n+1)-bakery runs = $_{def}$  (n)-bakery runs • bakery run

...

bakery run =<sub>def</sub> bake • supply to aide • move to shop • sell

(1)-bakery runs	= <sub>def</sub>	bakery run
(2)-bakery runs	= <sub>def</sub>	(1)-bakery runs • bakery run
(3)-bakery runs	= <sub>def</sub>	(2)-bakery runs • bakery run
•••		

(n+1)-bakery runs = $_{def}$  (n)-bakery runs • bakery run

This is an infinite process – that's really large...

#### A more realistic large process: order fulfillment



#### Some details



#### Some details



#### Some details


### Some details



### A more realistic large process: order fulfillment



### A more realistic large process: order fulfillment



just write A • B • C • D • E • F • G • H • I • J • K • L • M • N • O • P • Q • R • S • T • U • V • W • X • Y • Z

#### three challenges

- 1. objects and data
- 2. interaction between cases
- 3. composition

### c large process: order fulfillment



just write A • B • C • D • E • F • G • H • I • J • K • L • M • N • O • P • Q • R • S • T • U • V • W • X • Y • Z





# Excursus: Why is it important to distinguish between items ("real-world objects") and data?



#### Wolfgang can eat the *fish head*. But he cannot eat the *price* of the fish head!

price of a fish head ("data")

a fish head ("item")

43







### Agenda

# Large Language Models Process Models Large Process Models

Deutsches Forschungszentrum für Künstliche Intelligenz German Research Center for Artificial Intelligence



d

UNIVERSITÄT DES SAARLANDES

HUMBOLDT-UNIVERSITÄT ZU BERLIN



### **Commercial break**

### **Commercial break**

Peter Fettke · Wolfgang Reisig Understanding the Digital World Modeling with HERAKLIT

This book fills a serious gap by providing a conceptual framework for understanding the digital world. This world contains large, heterogeneous systems that have to manage dynamic behavior as well as static items and data. Obviously, new, *digital methods* are needed to deal with the challenges of the digital world.

This book introduces such a method with HERAKLIT, an intuitively simple, albeit powerful framework for modeling, communicating, and analyzing computerintegrated systems. It integrates proven methods for composing modules, describing behavior with local cause and effect, and digitally representing real- and imaginedworld items, resulting in a comprehensive, expressive, concerted, technically simple, digital modeling method.

This book is structured according to three HERAKLIT pillars, starting in Part I with the central HERAKLIT concept of modules, in particular their composition and refinement. Part II covers the second pillar of HERAKLIT, dynamics, focusing on modules that describe aspects of behavior. Part III focuses on static aspects. In particular, real- and imagined-world items and their symbolic representation are carefully distinguished and related. Together, these three pillars are consolidated in Part IV, integrating all concepts into a powerful formal framework. The book concludes in Part V with a more comprehensive case study of a typical retail business, recommendations on how to start modeling with HERAKLIT, and useful graphical conventions for the graphical representation of HERAKLIT, models.

HERAKLIT covers the range from the first informal structuring ideas for a computerintegrated system, through the specification of (business) processes, the contributions of people, organizations, and mechanical devices, up to the construction of software. The book is therefore written for students in areas related to system modeling, system design, and system engineering, as well as for professionals in these fields. Understanding the Digital World

Fettke · Reisig

Peter Fettke Wolfgang Reisig

### Understanding the Digital World

Modeling with HERAKLIT





### Idea for a large process model

idea of a large language model is applied in the context of business process management

### Idea for a large process model

idea of a large language model is applied in the context of business process management

three approaches:

- (1) fined-tuned large language model
- (2) large language model agent
- (3) native large language model



Enterprise Modelling and Information Systems Architectures Vol. 18, No. 3 (2023). DOI:10.18417/emisa.18.3 Editorial

Conceptual Modeling and Large Language Models: Impressions From First Experiments With ChatGPT

Hans-Georg Fill\*,a, Peter Fettkeb, Julius Köpkec

<sup>a</sup> University of Fribourg, Switzerland <sup>b</sup> Saarland University and German Research Center for Artificial Intelligence (DFKI), Germany <sup>c</sup> University of Klagenfurt, Austria

#### 1 Motivation

Since OpenAI publicly released ChatGPT in November 2022<sup>1</sup>, many ideas have emerged as to which applications this type of technology could support. At its core, ChatGPT is a *conversational artificial intelligence*, meaning that it can engage in a dialogue to respond to user input given in natural language (Campbell 2020). Although such types of systems have been well-known since Weizenbaum's Eliza program (Weizenbaum 1966) and are today widely deployed in practice under the popular term *chathots*, ChatGPT has a particular set of properties that contributed to its wide reception and the recent hype surrounding it.

In contrast to previous chatbots, ChatGPT does not retrieve responses from a knowledge base, which has been pre-defined by some human user. Rather, it is based on a pre-trained generative language model, which creates responses based on patterns that the user supplies as input. Thereby, a language model basically assigns probabilities to every word in a vocabulary that can follow a given input sequence. Such word embeddings are trained using artificial neural networks to learn a probability distribution from given texts in an unsupervised fashion, i.e. such that no additional human input or labeling is required. The generation of the output sequence thereby considers the tokens of the input sequence and their position as well as the previously generated output,

\* Corresponding author. E-mail. hans-georg.fill@unifr.ch 1 https://openai.com/blog/chatgpt which is thus denoted as an autoregressive generation (Juarfsky and Martin 2023). For the training of these probability distributions for the positional word embeddings, large sets of training data are required. In the case of the GPT-3 model, which underlies ChatGPT, this amounted to 175 billion parameters (Brown et al. 2020). For efficiently handling such large parameter sets, several innovations such as the architecture of transformer models (Vaswani et al. 2017) were necessary. 1

What seems to make ChatGPT however outpetform large language model (LLM) based programs that had been released to the public previously, is its ability to reduce toxic outputs, i. e. harmful or biased results. This has been achieved through the approach of InstructGPT (Ouyang et al. 2022), which uses reinforcement learning from human feedback to train a reward model. This reward model is then in turn used to fine-tune the output generated by the GPT-3 and GPT-4 language models. Thus, the training of the reward model only requires rather limited resources compared to the size of the language model.

From the multitude of areas in which the application of ChatGPT is currently discussed the following shall serve as non-exhaustive examples, which have already appeared in academic outlets: In the medical domain, ChatGPT has been considered for writing patient clinic letters by giving it instructions for following specific directions, using national guidelines and data from these guidelines in order to derive clinical advice (Ali et al. 2023); in the legal domain, ChatGPT has been asked to

Enterprise Modelling and Information Systems Architectures
Vol. 18, No. 3 (2023). DOI:10.18417/emisa.18.3
Editorial

Conceptual Modeling and Large Language Models: Impressions From First Experiments With ChatGPT

Hans-Georg Fill\*,a, Peter Fettkeb, Julius Köpkec

<sup>4</sup> University of Fribourg, Switzerland <sup>b</sup> Saarland University and German Research Center for Artificial Intelligence (DFKI), Germany <sup>c</sup> University of Klagnefurt, Austria

#### 1 Motivation

Since OpenAI publicly released ChatGPT in November 2022<sup>1</sup>, many ideas have emerged as to which applications this type of technology could support. At its core, ChatGPT is a *conversational artificial intelligence*, meaning that it can engage in a dialogue to respond to user input given in natural language (Campbell 2020). Although such types of systems have been well-known since Weizenbaum's Eliza program (Weizenbaum 1966) and are today widely deployed in practice under the popular term *chathots*, ChatGPT has a particular set of properties that contributed to its wide reception and the recent hype surrounding it.

In contrast to previous chatbots, ChatGPT does not retrieve responses from a knowledge base, which has been pre-defined by some human user. Rather, it is based on a pre-trained generative language model, which creates responses based on patterns that the user supplies as input. Thereby, a language model basically assigns probabilities to every word in a vocabulary that can follow a given input sequence. Such word embeddings are trained using artificial neural networks to learn a probability distribution from given texts in an unsupervised fashion, i.e. such that no additional human input or labeling is required. The generation of the output sequence thereby considers the tokens of the input sequence and their position as well as the previously generated output,

\* Corresponding author. E-mail. hans-georg.fill@unifr.ch <sup>1</sup> https://openai.com/blog/chatgpt which is thus denoted as an autoregressive generation (Jurafsky and Martin 2023). For the training of these probability distributions for the positional word embeddings, large sets of training data are required. In the case of the GPT-3 model, which underlies ChatGPT, this amounted to 175 billion parameters (Brown et al. 2020). For efficiently handling such large parameter sets, several innovations such as the architecture of transformer models (Vaswani et al. 2017) were necessary. 1

What seems to make ChatGPT however outperform large language model (LLM) based programs that had been released to the public previously, is its ability to reduce toxic outputs, i. e. harmful or biased results. This has been achieved through the approach of InstructGPT (Ouyang et al. 2022), which uses reinforcement learning from human feedback to train a reward model. This reward model is then in turn used to fine-tune the output generated by the GPT-3 and GPT-4 language models. Thus, the training of the reward model only requires rather limited resources compared to the size of the language model.

From the multitude of areas in which the application of ChatGPT is currently discussed, the following shall serve as non-exhaustive examples, which have already appeared in academic outlets: In the medical domain, ChatGPT has been considered for writing patient clinic letters by giving it instructions for following specific directions, using national guidelines and data from these guidelines in order to derive clinical advice (Ali et al. 2023); in the legal domain, ChatGPT has been asked to

#### **Task Definition 5**

Lighting and ventilation of a bathroom: If the light is switched on when the fan is stationary, the fan also starts after a while. Then, when the light is turned off, the fan continues to run for some time. If the light is turned on first and then turned off quickly when the fan is stationary, the fan will not start at all. If the light is first switched off and then quickly switched on again when the fan is running, the fan continues to run without interruption.









large language model

large language model agent

large language model













*(On the official website of an airline)* Book the cheapest flight from Beijing to Los Angeles in the last week of July.



(In the middle of a kitchen in a simulator) Please put a pan on the dinning table.

*(On the official website of an airline)* Book the cheapest flight from Beijing to Los Angeles in the last week of July.





#### **Real-world Challenges**

(On an Ubuntu bash terminal) Recursively set all files in the directory to read-only, except those of mine.

(Given Freebase APIs) What musical instruments do Minnesotaborn Nobel Prize winners play?

(Given MySQL APIs and existed tables) Grade students over 60 as PASS in the table.

(On the GUI of Aquawar) This is a two-player battle game, you are a player with four pet fish cards .....

A man walked into a restaurant, ordered a bowl of turtle soup, and after finishing it, he committed suicide. Why did he do that?

(In the middle of a kitchen in a simulator) Please put a pan on the dinning table.

*(On the official website of an airline)* Book the cheapest flight from Beijing to Los Angeles in the last week of July.

#### 8 Distinct Environments



#### **Real-world Challenges**



#### **8** Distinct Environments

### (3) Large process model as a native large language model



### (3) Large process model as a native large language model


### (3) Large process model as a native large language model

	Contents lists available at ScienceDirect	The latent for the second seco		
- CL	Decision Support Systems			
EL CEVIED	journal homepage: www.elsevier.com/locate/dss			
Predicting proce	sss behaviour using deep learning	Cros		
Predicting proce	rss behaviour using deep learning ana-Rebecca Rehse <sup>b, c</sup> , Peter Fettke <sup>b, c</sup>	Cros		
Predicting proce Joerg Evermann <sup>a,*</sup> , J. <sup>4</sup> Ammial Indiversity of Newford <sup>6</sup> Gaman Research Center for Artific <sup>6</sup> Saarbridsen,	ess behaviour using deep learning ana-Rebecca Rehse <sup>b, c</sup> , Peter Fettke <sup>b, c</sup> land S. Johrs M. Canada ici Intelligence, Saarbrücken, Germany Germany	Cros		

state-of-the-art in prediction precision.

Keywords: Process management Runtime support Process prediction Deep learning Neural networks

### 1. Introduction

Being able to predict the future behaviour of a business process is an important business capability [1]. As an application of predictive analytics in business process management, process prediction exploits data on past process instances to make predictions about current ones [2]. Example use cases are customer service agents responding to inquiries about the remaining time until a case is resolved, production managers predicting the completion time of a production process for better planning and higher utilization, or case managers identifying likely compliance violations to mitigate business risk.

We present a novel approach to predicting the next process event using deep learning. While the term "deep learning" has only recently become a popular research topic, it is essentially an application of neural networks and thus looks back on a long history of research [3]. Recent innovations both in algorithms, allowing novel architectures of neural networks, and computing hardware, especially GPU processing, have led to a resurgence in interest for neural networks and popularized the term "deep learning" [4]. Our approach is motivated by applications of neural networks to natural language processing(NLP), more specifically the prediction of the next word in a sentence [5-7]. By interpreting process event logs

\* Corresponding author. E-mail address: jevermann@mun.ca (]. Evermann).

http://dx.doi.org/10.1016/i.dss.2017.04.003 0167-9236/© 2017 Elsevier B.V. All rights reserved.

as text, process traces as sentences, and process events as words, these techniques can be applied to predict future process events. The contribution of our research is threefold:

1. We improve on the state-of-the-art in process event prediction. Our results show our method has considerably better

- sary for prediction. Deep learning models, where the process
- 3. We contribute to process management in general by showcasing the useful application of an artificial intelligence approach, from the application of smart approaches.

Our research is located at the intersection of business process management, in particular process mining, and artificial intelligence (AI) and machine learning. We bring together historic process data with an AI learning technology to leverage real-time case management, opening new perspectives into process execution, monitoring, and analysis. Extending existing solutions to novel problems ("exaptation") is a recognized and valid way to make a contribution in design science [8], which is the research approach we apply here. We not only provide a new approach, rooted in AI, to predicting the next

precision on next-event prediction. 2. We demonstrate that an explicit process model is not neces-

© 2017 Elsevier B.V. All rights reserved

- structure is only implicitly reflected, can perform as well as explicit process models.
- illustrating that business process management can benefit

73

### (3) Large process model as a native large language model

Contrato lista susilable et Spisson Disert	

Decision Support Systems

nar nomepage: www.eisevier.com/ioi

ecision Support Systems 100 (2017) 129–140



Joerg Evermann<sup>a,\*</sup>, Jana-Rebecca Rehse<sup>b, c</sup>, Peter Fettke<sup>b, c</sup>

<sup>a</sup> Memorial University of Newfoundland, St. John's, NL, Canada
<sup>b</sup> German Research Center for Artificial Intelligence, Saarbrücken, Germany
<sup>c</sup> Saarland University, Saarbrücken, Germany

### ARTICLE INFO ABSTRACT

Article history: Received 8 July 2016 Received in revised form 22 March 2017 Accepted 5 April 2017 Available online 17 April 2017

Keywords: Process management Runtime support Process prediction Deep learning Neural networks Predicting business process behaviour is an important aspect of business process management. Motivated pur research in narrul alnguage processing, this paper describes an application of deep learning with recurrent neural networks to the problem of predicting the next event in a business process. This is both a novel method in process prediction, which has largely relead on explicit process models, and also a novel application of deep learning methods. The approach is evaluated on two real datasets and our results surpass the state-of-the-art in mediction process.

© 2017 Elsevier B.V. All rights reserved.

CrossMar

### 1. Introduction

Being able to predict the future behaviour of a business process is an important business capability [1]. As an application of predictive analytics in business process management, process prediction exploits data on past process instances to make predictions about current ones [2]. Example use cases are customer service agents responding to inquiries about the remaining time until a case is resolved, production managers predicting the completion time of a production process for better planning and higher utilization, or case managers identifying likely compliance violations to mitigate business risk.

We present a novel approach to predicting the next process event using deep learning. While the term 'deep learning' has only recently become a popular research topic, it is essentially an application of neural networks and thus looks back on a long history of research [3]. Recent innovations both in algorithms, allowing novel architectures of neural networks, and computing hardware, especially CPU processing, have led to a resurgence in interest for neural networks and popularized the term 'deep learning' [4]. Our approach is motivated by applications of neural networks to natural language processing(NLP), more specifically the prediction of the next word in a sentence [5–7]. By interpring process event logs

\* Corresponding author. E-mail address: jevermann@mun.ca (J. Evermann).

http://dx.doi.org/10.1016/j.dss.2017.04.003 0167-9236/© 2017 Elsevier B.V. All rights reserved. as text, process traces as sentences, and process events as words, these techniques can be applied to predict future process events. The contribution of our research is threefold:

 We improve on the state-of-the-art in process event prediction. Our results show our method has considerably better precision on next-event prediction.
 We demonstrate that an explicit process model is not neces-

- sary for prediction. Deep learning models, where the process structure is only implicitly reflected, can perform as well as explicit process models.
- We contribute to process management in general by showcasing the useful application of an artificial intelligence approach, illustrating that business process management can benefit from the application of smart approaches.

Our research is located at the intersection of business process management, in particular process mining, and artificial intelligence (AI) and machine learning. We bring together historic process data with an AI learning technology to leverage real-line case management, opening new perspectives into process execution, monitoring, and analysis. Extending existing solutions to novel problems ("exaptation") is a recognized and valid way to make a contribution in design science [8], which is the research approach we apply here. We not only provide a new approach, rooted in AI, to predicting the next Artificial Intelligence Review (2022) 55:801–827 https://doi.org/10.1007/s10462-021-09960-8



A systematic literature review on state-of-the-art deep learning methods for process prediction

Dominic A. Neu<sup>1,2</sup> · Johannes Lahann<sup>1,2</sup> · Peter Fettke<sup>1,2</sup>

Published online: 11 March 2021 © The Author(s) 2021

### Abstract

Process mining enables the reconstruction and evaluation of business processes based on digital traces in IT systems. An increasingly important technique in this context is process prediction. Given a sequence of events of an ongoing trace, process prediction allows forecasting upcoming events or performance measurements. In recent years, multiple process prediction approaches have been proposed, applying different data processing schemes and prediction algorithms. This study focuses on deep learning algorithms since they seem to outperform their machine learning alternatives consistently. Whilst having a common learning algorithm, they use different data preprocessing techniques, implement a variety of network topologies and focus on various goals such as outcome prediction, time prediction or control-flow prediction. Additionally, the set of log-data, evaluation metrics and baselines used by the authors diverge, making the results hard to compare. This paper attempts to synthesise the advantages and disadvantages of the procedural decisions in these approaches by conducting a systematic literature review.

Keywords Process prediction  $\cdot$  Predictive process monitoring  $\cdot$  Systematic literature review  $\cdot$  Deep learning

### 1 Introduction

Today's information systems create, utilize and store vast amounts of data about the business processes being executed with them. These logs capture the as-is execution. Process mining extracts knowledge from these logs to provide means for process discovery, process monitoring and process improvement. Additionally, in case target process models are provided, conformance-checking searches for deviations of this model. Consequently, process mining is situated between the disciplines of data mining and business process modelling (Van der Aalst et al. 2011). In recent years, much effort was put into process discovery to build human-readable models for further investigation by

dominic.neu@dfki.de

<sup>🖂</sup> Dominic A. Neu

<sup>&</sup>lt;sup>1</sup> Institute for Information Systems, German Research Center for Artificial Intelligence (DFKI), Saarbruecken, Germany

<sup>&</sup>lt;sup>2</sup> Institute for Information Systems, Saarland University, Saarbruecken, Germany

### (3) Large process model as a native large language model

### Decision Support Systems 100 (2017) 129-140

Contents lists available at ScienceDirect Decision Support Systems journal homepage: www.elsevier.com/locate/dss

state-of-the-art in prediction precision.

Predicting process behaviour using deep learning

Joerg Evermann<sup>a,\*</sup>, Jana-Rebecca Rehse<sup>b, c</sup>, Peter Fettke<sup>b, c</sup>

a Memorial University of Newfoundland, St. John's, NL, Canada b German Research Center for Artificial Intelligence, Saarbrücken, Germany Saarland University, Saarbrücken, Germany

### ARTICLE INFO ABSTRACT

Article history Received 8 July 2016 Received in revised form 22 March 2017 Accepted 5 April 2017 Available online 17 April 2017

Kennwords: Process management Runtime support Process prediction Deep learning Neural networks

### 1. Introduction

Being able to predict the future behaviour of a business process contribution of our research is threefold is an important business capability [1]. As an application of predictive analytics in business process management, process prediction exploits data on past process instances to make predictions about current ones [2] Example use cases are customer service agents responding to inquiries about the remaining time until a case is resolved, production managers predicting the completion time of a production process for better planning and higher utilization, or case managers identifying likely compliance violations to mitigate business risk.

We present a novel approach to predicting the next process event using deep learning. While the term "deep learning" has only recently become a popular research topic, it is essentially an application of neural networks and thus looks back on a long history of research [3]. Recent innovations both in algorithms, allowing novel architectures of neural networks, and computing hardware, especially GPU processing, have led to a resurgence in interest for neural networks and popularized the term "deep learning" [4]. Our approach is motivated by applications of neural networks to natural language processing(NLP), more specifically the prediction of the next word in a sentence [5-7]. By interpreting process event logs

\* Corresponding author. E-mail address: jevermann@mun.ca (]. Evermann).

http://dx.doi.org/10.1016/i.dss.2017.04.003 0167-9236/© 2017 Elsevier B.V. All rights reserved as text, process traces as sentences, and process events as words, these techniques can be applied to predict future process events. The

© 2017 Elsevier B.V. All rights reserved.

Predicting business process behaviour is an important aspect of business process management. Motivated

rent neural networks to the problem of predicting the next event in a business process. This is both a novel

method in process prediction, which has largely relied on explicit process models, and also a novel applica-

tion of deep learning methods. The approach is evaluated on two real datasets and our results surpass the

by research in natural language processing, this paper describes an application of deep learning with recur-

CrossMark

1. We improve on the state-of-the-art in process event prediction. Our results show our method has considerably better precision on next-event prediction 2. We demonstrate that an explicit process model is not neces-

- sary for prediction. Deep learning models, where the process structure is only implicitly reflected, can perform as well as explicit process models
- 3. We contribute to process management in general by showcasing the useful application of an artificial intelligence approach. illustrating that business process management can benefit from the application of smart approaches.

Our research is located at the intersection of business process management, in particular process mining, and artificial intelligence (AI) and machine learning. We bring together historic process data with an AI learning technology to leverage real-time case management, opening new perspectives into process execution, monitoring, and analysis. Extending existing solutions to novel problems ("exaptation") is a recognized and valid way to make a contribution in design science [8], which is the research approach we apply here. We not only provide a new approach, rooted in AI, to predicting the next Artificial Intelligence Review (2022) 55:801-827 https://doi.org/10.1007/s10462-021-09960-8

A systematic literature review on state-of-the-art deep learning methods for process prediction

Dominic A. Neu<sup>1,2</sup> · Johannes Lahann<sup>1,2</sup> · Peter Fettke<sup>1,2</sup>

Published online: 11 March 2021 © The Author(s) 2021

### Abstract

Process mining enables the reconstruction and evaluation of business processes based on digital traces in IT systems. An increasingly important technique in this context is process prediction. Given a sequence of events of an ongoing trace, process prediction allows forecasting upcoming events or performance measurements. In recent years, multiple process prediction approaches have been proposed, applying different data processing schemes and prediction algorithms. This study focuses on deep learning algorithms since they seem to outperform their machine learning alternatives consistently. Whilst having a common learning algorithm, they use different data preprocessing techniques, implement a variety of network topologies and focus on various goals such as outcome prediction, time prediction or control-flow prediction. Additionally, the set of log-data, evaluation metrics and baselines used by the authors diverge, making the results hard to compare. This paper attempts to synthesise the advantages and disadvantages of the procedural decisions in these approaches by conducting a systematic literature review.

Keywords Process prediction · Predictive process monitoring · Systematic literature review · Deep learning

### 1 Introduction

Today's information systems create, utilize and store vast amounts of data about the business processes being executed with them. These logs capture the as-is execution. Process mining extracts knowledge from these logs to provide means for process discovery, process monitoring and process improvement. Additionally, in case target process models are provided, conformance-checking searches for deviations of this model. Consequently, process mining is situated between the disciplines of data mining and business process modelling (Van der Aalst et al. 2011). In recent years, much effort was put into process discovery to build human-readable models for further investigation by

dominic.neu@dfki.de

Institute for Information Systems, German Research Center for Artificial Intelligence (DFKI), Saarbruecken, Germany

<sup>2</sup> Institute for Information Systems, Saarland University, Saarbruecken, Germany



### **PGTNet: A Process Graph Transformer Network** for Remaining Time Prediction of Business **Process Instances**

Keyvan Amiri Elyasi<sup>1[0009-0007-3016-2392]</sup>, Han van der Aa<sup>2</sup>[0000-0002-4200-4937], and Heiner Stuckenschmidt<sup>1</sup>[0000-0002-0209-3859]

<sup>1</sup> Data and Web Science Group, University of Mannheim, Germany {keyvan,heiner}@informatik.uni-mannheim.de <sup>2</sup> Faculty of Computer Science, University of Vienna, Austria han.van.der.aa@univie.ac.at

Abstract. We present PGTNet, an approach that transforms event logs into graph datasets and leverages graph-oriented data for training Process Graph Transformer Networks to predict the remaining time of business process instances. PGTNet consistently outperforms state-of-the-art deep learning approaches across a diverse range of 20 publicly available real-world event logs. Notably, our approach is most promising for highly complex processes, where existing deep learning approaches encounter difficulties stemming from their limited ability to learn control-flow relationships among process activities and capture long-range dependencies. PGTNet addresses these challenges, while also being able to consider multiple process perspectives during the learning process.

Keywords: Predictive process monitoring · Remaining time prediction · Deep learning · Graph Transformers.

### 1 Introduction

Predictive process monitoring (PPM) aims to forecast the future behaviour of running business process instances, thereby enabling organizations to optimize their resource allocation and planning [17], as well as take corrective actions [7]. An important task in PPM is remaining time prediction, which strives to accurately predict the time until an active process instance will be completed. Precise estimations for remaining time are crucial for avoiding deadline violations, optimizing operational efficiency, and providing estimates to customers [13, 17].

A variety of approaches have been developed to tackle remaining time prediction, with recent works primarily being based on deep learning architectures. In this regard, approaches using deep neural networks are among the most prominent ones [15]. However, the predictive accuracy of these networks leaves considerable room for improvement. In particular, they face challenges when it comes to capturing long-range dependencies [2] and other control-flow relationships (such as loops and parallelism) between process activities [22], whereas they

Dominic A Neu

### Table 2. Mean Absolute Error for remaining time prediction (MAE: in days).

Decision Support Systems IT Contents lists available Decision Support ELSEVIER journal homepage: www.el

(3) Large

Predicting process behaviour using deep lear

ABSTRACT

Predicting business process beh by research in natural language

rent neural networks to the prob

method in process prediction,

tion of deep learning methods.

state-of-the-art in prediction pre-

Joerg Evermann<sup>a</sup>.\*, Jana-Rebecca Rehse<sup>b, c</sup>, Peter Fettke<sup>b</sup>
<sup>a</sup> Memorial University of Newfoundland, St. John's, NL. Canada
<sup>b</sup> German Research Center for Artificial Intelligence, Saubriticken, Germany
<sup>c</sup> Saudnad University, Saubriticken, Germany

### ARTICLE INFO

Article history: Received 19 July 2016 Received in revised form 22 March 2017 Accepted 5 April 2017 Available online 17 April 2017

Keywords: Process management Runtime support Process prediction Deep learning Neural networks

### 1. Introduction

Being able to predict the future behaviour of a business process is an important business capability [1]. As an application of predictive analytics in business process management, process prediction exploits data on past process instances to make predictions about current ones [2]. Example use cases are customer service agents resolved, production managers predicting the completion time of a production process for better planning and higher utilization, or case managers identifying likely compliance violations to mitigate business risk.

We present a novel approach to predicting the next process event using deep learning. While the term 'deep learning' has only recently become a popular research topic, it is essentially an application of neural networks and thus looks back on a long history of research [3]. Recent innovations both in algorithms, allowing novel architectures of neural networks, and computing hardware, especially GPU processing, have led to a resurgence in interest for neural networks and popularized the term 'deep learning' [4]. Our approach is motivated by applications of neural networks to natural language processing(NLP), more specifically the prediction of the next word in a sentence [5–7]. By interpreting process vent logs

\* Corresponding author. E-mail address: jevermann@mun.ca (J. Evermann).

http://dx.doi.org/10.1016/j.dss.2017.04.003 0167-9236/© 2017 Elsevier B.V. All rights reserved.

Event log	DIIMMV	DALSTM	Process	CCNN	PGTNet	
		DALSIM	Transformer	GGINI		
Env.permit	5.21	$3.36\pm~0.04$	$4.26 \pm \ 0.04$	$3.52\pm~0.02$	$2.72\pm$	0.08
$\operatorname{Helpdesk}$	9.15	$8.22{\pm}~0.23$	$6.33\pm~0.01$	$6.21\pm~0.04$	$4.11\pm$	0 <b>.04</b>
BPIC12	9.03	$9.34 \pm \ 0.41$	$7.11\pm~0.02$	$4.78{\pm}~0.01$	$2.31 \pm$	0 <b>.19</b>
BPIC12W	9.16	$8.22{\pm}~0.06$	$7.40\pm~0.01$	$5.12\pm~0.02$	$2.70\pm$	0 <b>.01</b>
BPIC12C	8.92	$8.21{\pm}~0.27$	$6.86\pm~0.01$	$5.32\pm~0.01$	$2.77\pm$	0 <b>.02</b>
BPIC12CW	9.17	$8.04{\pm}~0.09$	$7.46 \pm \hspace{0.1cm} 0.01$	$6.99\pm~0.01$	$5.07\pm$	0 <b>.03</b>
BPIC12O	8.39	$8.21{\pm}~0.09$	$7.29 \pm \ 0.01$	$6.93 \pm 0.04$	$5.57\pm$	0 <b>.01</b>
BPIC12A	8.17	$7.62{\pm}~0.03$	$7.79\pm~0.01$	$7.48\pm~0.01$	$7.38\pm$	0 <b>.01</b>
BPIC20I	27.20	$20.43 \pm \hspace{0.1cm} 0.39$	$17.06 \pm 0.11$	$15.67 \pm 0.04$	$7.67\pm$	0 <b>.19</b>
BPIC20D	4.33	$4.15{\pm}~0.12$	$3.65\pm~0.01$	$3.25\pm~0.01$	$3.10\pm$	0 <b>.01</b>
$\operatorname{Sepsis}$	41.12	$25.21{\pm}~0.66$	$34.77 \pm 0.18$	$19.44 \pm \ 0.05$	$16.48 \pm$	0 <b>.19</b>
$\operatorname{Hospital}$	59.41	$43.66 \pm 0.10$	$47.00 \pm 0.07$	$41.84 \pm 0.06$	$35.68 \pm$	0 <b>.03</b>
BPIC15-1	50.22	$36.48\pm~2.69$	$31.01{\pm}~0.36$	$16.77 \pm 0.01$	$1.76\pm$	0 <b>.06</b>
BPIC15-2	83.11	$63.66 \pm \ 2.36$	$44.04{\pm}~0.48$	$20.76 \pm \ 0.05$	$3.02\pm$	0 <b>.07</b>
BPIC15-3	28.76	$17.69 \pm 1.16$	$15.23{\pm}~0.23$	$7.06\pm~0.03$	$1.54\pm$	0.23
BPIC15-4	56.75	$53.33 \pm \ 2.63$	$34.40{\pm}~0.42$	$17.97 \pm 0.03$	$1.65\pm$	0 <b>.06</b>
BPIC15-5	45.97	$42.89 \pm \hspace{0.1cm} 3.08$	$27.76 \pm 0.28$	$13.61 \pm 0.08$	$1.61\pm$	0 <b>.01</b>
BPIC13I	16.18	$7.60{\pm}~0.45$	$13.54\pm~0.04$	$11.99 \pm 0.04$	$2.23\pm$	0.05
BPIC13C	152.93	$91.82{\pm}~1.48$	$127.01 \pm 0.85$	$123.28 \pm 0.53$	$37.44 \pm$	1 <b>.49</b>
Traffic fines	196.26	$187.41 \pm 0.53$	$187.08 \pm 0.11$	$154.56 \pm 0.19$	$113.53\pm$	0 <b>.12</b>
Average	41.47	$32.78 \pm 0.84$	$31.85 \pm 0.16$	$24.63 \pm 0.06$	$12.92\pm$	0.14

model

Transformer Network diction of Business tances

-3016-2392], Han van der ekenschmidt<sup>1</sup>[0000-0002-0209-3859]

ersity of Mannheim, Germany k.uni-mannheim.de iversity of Vienna, Austria ivie.ac.at

roach that transforms event logs -oriented data for training Prodict the remaining time of busitly outperforms state-of-the-art e range of 20 publicly available ach is most promising for highly learning approaches encounter bility to learn control-flow relaupture long-range dependencies. ile also being able to consider learning process.

ng · Remaining time prediction

to forecast the future behaviour of enabling organizations to optimize us well as take corrective actions [7]. In prediction, which strives to accus instance will be completed. Precise r avoiding deadline violations, optiestimates to customers [13, 17]. ped to tackle remaining time predicd on deep learning architectures. In etworks are among the most promiacy of these networks leaves considthey face challenges when it comes and other control-flow relationships rocess activities [22], whereas they

### Central theorem (based on [Fettke / Reisig 2024])

Each finite process R can be composed as  $R = P_1 \bullet \dots \bullet P_n$  from steps  $P_1, \dots, P_n$ .

Steps  $P_1$ , ...,  $P_n$  can be interpreted as *tokens* ("the vocabulary of a system") and conventional methods for large language models can be employed.

### Conclusions

(A) large language models evolve into a foundational technology

## Conclusions

(A) large language models evolve into a foundational technology

(B) three approaches for large process models:

- (1) fine-tuned large language models
- (2) large language model agents
- (3) native large process models

## Conclusions

(A) large language models evolve into a foundational technology

(B) three approaches for large process models:

- (1) fine-tuned large language models
- (2) large language model agents
- (3) native large process models

(C) composition calculus is the theoretical foundation for large process models



# For a small obolus to Wolfgang, I will gladly answer your question!

### **Professor Dr. Peter Fettke**

Saarland University and German Research Center for Artificial Intelligence (DFKI) Campus D 3 2 66123 Saarbrücken, Germany peter.fettke@dfki.de, http://bpm.dfki.de Phone: +49 681 85775-5142

UNIVERSITÄT DES SAARLANDES

**Deutsches** 

Intelligenz

Intelligence

für Künstliche

Forschungszentrum

German Research Center for Artificial