# Evaluating the Process Modeling Abilities of Large Language Models – Preliminary Foundations and Results
## Research in Progress

Peter Fettke and Constantin Houy

German Research Center for Artificial Intelligence (DFKI) and
Saarland University, Saarbrücken, Germany
{peter.fettke,constantin.houy}@dfki.de

**Abstract.** Large language models (LLM) have revolutionized the processing of natural language. Although first benchmarks of the process modeling abilities of LLM are promising, it is currently under debate to what extent an LLM can generate good process models. In this contribution, we argue that the evaluation of the process modeling abilities of LLM is far from being trivial. Hence, available evaluation results must be taken carefully. For example, even in a simple scenario, not only the quality of a model should be taken into account, but also the costs and time needed for generation. Thus, an LLM does not generate one optimal solution, but a set of Pareto-optimal variants. Moreover, there are several further challenges which have to be taken into account, e.g. conceptualization of quality, validation of results, generalizability, and data leakage. We discuss these challenges in detail and discuss future experiments to tackle these challenges scientifically.

**Keywords:** large process models, model generation, automated modeling, BPMN, Pareto front

## 1 Motivation

Large language models (LLM) have revolutionized many tasks of natural language processing (NLP). Recently, LLM are not only used for language related tasks, but also for generating program code (Joel et al. 2024, Sarkar et al. 2022), planning (Katz et al. 2024, Erdogan et al. 2025), action models (Wang et al. 2025, Zhang et al. 2024), and many more.

In 2023, first ideas to use LLM for generating and understanding conceptual models emerged (Fill et al. 2023). Since then, different approaches for the use of LLM have been developed and tested. The latest developments aim at using LLM for interactive process modeling, e.g. (Kourani, Berti, Schuster & van der Aalst 2024*b*), for process discovery from event logs (Norouzifar et al. 2024), for improving process model understandability (Kourani, Berti, Hennrich, Kratsch, Weidlich, Li, Arslan, Schuster & van der Aalst 2024), and other tasks related to process modeling (Vidgof et al. 2023).

Although the machine-generated results are often surprising and of astonishing quality, it is obvious that the evaluation has to be conducted in some objective way. For such a systematic evaluation, already first benchmark studies have been undertaken and published (Kourani, Berti, Schuster & van der Aalst 2024*a*).

However, in this article, we argue that an objective evaluation is far from trivial for a number of reasons, e.g. model quality is multi-criteria measure (Gutschmidt & Nast 2025, Krogstie 2016, Heggset et al. 2015) with dimensions such as syntactic correctness, semantic adequacy, and understandability (Houy et al. 2012, Pavlicek et al. 2019). The long-term objective of this piece of research in progress aims at developing the necessary foundations and methods for an objective evaluation of the process modeling abilities of LLM. As a first step to tackle the evaluation challenges, we define a typical standard evaluation scenario. Based on this standard scenario, we are able to elaborate on particular evaluation challenges and possible ways to overcome them.

In the following, we report some preliminary results, namely, the assumptions of the standard evaluation scenario (section 2); the conceptualization of quality measures describing the LLM performance and the trade-off between quality, cost, and time (section 3); a specific evaluation example (section 4); the discussion of related work (section 5); and conclusions and discussion of future work (section 6).

## 2 Assumptions of the standard evaluation scenario

We introduce three assumptions to define the standard evaluation scenario:

1. A process model conceptualizes a real or imagined modeling domain. We assume that the domain to be modeled is given by a plain English text. In other words, the task to be accomplished by the LLM is to transform or reformulate the natural textual description by using a process modeling language.
2. For each textual domain description one or more sample solutions are given. They serve as a basis for model evaluation and can be considered as the "gold-standard"; in machine learning, the gold-standard is typically called "ground-truth".
3. The sample solutions of the gold-standard use the Business Process Modeling Notation (BPMN) as a process modeling language.

It is obvious that in many real-world modeling situations the three assumptions are not fulfilled. However, many arguments justify them, e.g. BPMN is often treated as the de-facto standard for process modeling (von Rosing et al. 2015), a textual domain description provides a controlled environment which is necessary for the comparison of results, numerous examples demonstrate that the generation of a process model based on a textual description is realistic to a certain degree and feasible for humans and machines. However, later we will discuss some consequences for the evaluation if one or more of these assumptions are relaxed in particular ways.

## 3 Measuring the performance in the standard evaluation scenario

### 3.1 Quality

Results from the field of NLP clearly demonstrate that measuring the quality of a machine-generated language artifact and comparing it with the quality of a human-generated task is far from simple, e.g. evaluating the subtle language differences in

machine translation. We argue that the challenges of evaluation in the standard process modeling scenario are comparable to the challenges known in the field of NLP.

One simple solution, also often used in NLP research, relies on the opinions of domain or modeling experts. Such experts are often acquired using well-known crowd-working platforms. Such an approach is inherently subjective and it is not clear how to control the quality of the experts' work, particular in large scale evaluations.

More objective measures are also problematic. While syntactic deficits can relatively easily be counted, semantic differences are much more difficult to evaluate, e.g. different abstraction levels used for modeling, focusing on relevant aspects or the omission of irrelevant details. Last but not least, it could be argued that pragmatic aspects of the model must also be taken into account, e.g. the layout of a BPMN diagram, the intended purpose of the model or the task the modeller would like to solve with the model.

Similar to typical measures used in information retrieval, e.g. precison and recall, the quality of the model in question could be measured: How many concepts of the model are not contained in the gold-standard? How many concepts of the gold-standard are not in the model?

The previous idea can be conceptualized more concrete, if an execution model of a BPMN diagram is additionally assumed: It could be argued that all execution paths of the gold-standard should be included in the presented solution (Kourani, Berti, Schuster & van der Aalst 2024*a*). On the other hand, execution paths of a particular solution which are not possible in the gold-standard should be penalized. Note, that such an approach relies on the use of an appropriate matching between the gold-standard and the provided solution – and process matching is far from easy (Thaler et al. 2014).

Since BPMN diagrams can be interpreted as graphs, the idea of a graph-edit-distance (Bunke & Shearer 1998) can be employed (Dijkman et al. 2011): How many edit operations are necessary to transform the provided solution into the gold-standard model? This idea is quite simple and can theoretically be easily interpreted: the quality of a machine-generated solution is correlated to the amount of rework needed to revise the model accordingly. The particular needed edit operations could also be weighted based on a specific cost model for each operation.

Currently, although the described ideas for measuring the quality has some face validity, the validity of these measures for evaluating process modeling abilitites of LLM is more or less unknown.

## 3.2 Time and costs

In addition to the quality of the generated solution, the time for generating a model is of importance. For such analyses, theoretical computer science developed a rich theory supporting the analysis of algorithms based on different inputs; the so-called "Big O notation". However, currently we do not know if such a theoretical analysis produces valid results in our case and we are not aware of any work which is relevant in this context. As a simple alternative, the time an LLM needed for generating the output model can easily be measured in an evaluation procedure.

Furthermore, the usage of LLM is not without costs. First, there are some fix costs for setting up the infrastructure and training the LLM. Second, there are variable costs

for using LLM. One major idea of LLM are so-called foundational models which can be used for different tasks and are offered as services by different technology companies (Bommasani et al. 2022, Schneider et al. 2024). The variable costs for usage are typically measured in input tokens, output tokens, API calls, etc., and have to be paid in US-Dollars, Euros, or other currencies. Hence, it is also possible to quantify the costs of using an LLM in such service scenarios. However, note, that the particular variable costs of using foundational models in evaluating the abilities of LLM should not be neglected; e.g. Stein et al. (2025) do not benchmark every combination of an experimental study design for evaluating the planning abilities of LLM because the estimated costs exceeded their budget. This is not without practical relevance, particular in comparison for the labor costs needed for a human-generated model.

One further aspect is of importance: Since LLM do not work deterministically, a repetitive measurement of all measures is needed and some averaging is necessary. Therefore, systematic studies of LLM parameters such as the so-called "seed" and "temperature" of the model must be studied in a systematic way.

## 4 An example

So far, we have not established the full infrastructure for presenting a completed benchmarking study. Hence, the presentation of the results of a full automation of the process is currently not possible. However, we aim at a fully-automated pipeline for the above-mentioned standard case. Previous work on automated assessment of modeling results demonstrates that automatically measuring quality is possible in principle, e.g. Thaler et al. (2016), Ullrich et al. (2023). In addition to that, the cost to compute the benchmarks plays an important role. However, the development of approaches is progressing rapidly and many different fine-tuned LLM supporting different specific tasks exist. Against this background, it is currently questionable how science can keep up to produce a truly acceptable assessment of approaches. In turn, this does not mean that we should not try. Documented experience with several LLM shows that for some textual descriptions, acceptable or even good results were produced. Inspired by the study of Kourani, Berti, Schuster & van der Aalst (2024*a*), we used their results on the quality and time efficiency of "a diverse set of state-of-the-art LLM" for process modeling tasks [p.11] and extended them by adding the cost dimension measured in US-Dollars. In this context, we let GPT-4 estimate the token consumption of each LLM based on the given average time elapsed for producing an acceptable model in the above study, checked them for plausibility, and calculated the resulting costs (in US-Dollars) using an extensive and commonly-used LLM pricing calculator (https://huggingface.co/spaces/Presidentlin/llm-pricing-calculator) based on the prices of providers of each LLM. The following diagrams visualize our results in combination with the results of Kourani, Berti, Schuster & van der Aalst (2024*a*) and present three Pareto fronts for the perspectives "quality vs. cost", "quality vs. time" and "cost vs. time". In the latter, there is one optimal data point concerning cost vs. time, hence the Pareto front is degenerated and there is only one Pareto-optimal solution in this example (see Figure 1).
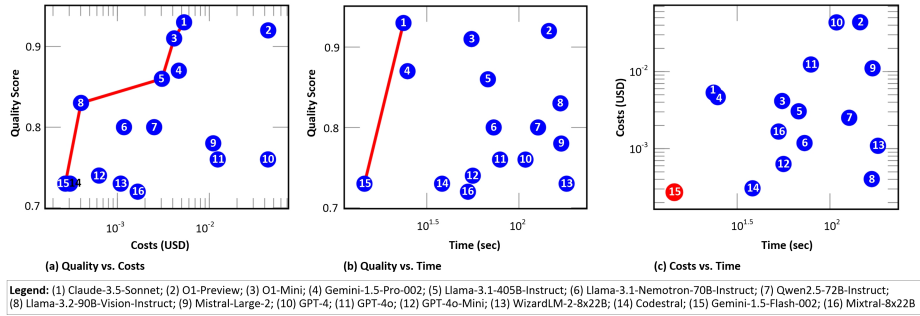
**Figure 1.** Pareto analysis for Quality, Cost, and Time (Pareto fronts in red, log-scaled axes)

## 5 Related work

We relate our work to five research streams, namely (1) automated modeling, (2) automated model assessment, (3) evaluation of abilities of LLM in general, (4) known evaluations of process modeling abilities of LLM, and (5) available data sets.

(1) While early work on automated modeling used rule-based approaches, recent articles investigate the potential of LLM for generating process models from the content of documents (Voelter et al. 2024) or using other formats of process description, like declarative process descriptions which can be automatically transformed into BPMN models (Wiśniewski et al. 2019). (2) An overview of works on the automated assessment of conceptual models as well as process models in education can be found in Ullrich et al. (2023). For example, Thaler et al. (2016) present an approach for the automated assessment of process model quality with regard to syntax, semantics, and pragmatic quality, e.g. understandability, for the example of Event-driven Process Chains (EPC). Westergaard et al. (2013) present an according approach for coloured Petri nets, Sanchez-Ferreres et al. (2020) for BPMN, and Schramm et al. (2012), Striewe & Goedicke (2014) and Beck et al. (2015) focus on UML activity diagrams in the context of teaching process modeling. (3) The evaluation of abilities of LLM in general is treated in many different works proposing different approaches for different domains, e.g. Hu & Zhou (2024), Grandi et al. (2024). Furthermore, evaluation foundations are generally discussed, e.g. by Mizrahi et al. (2024) and Kapoor et al. (2024). (4) The so far most pertinent work in the context of studying the evaluation of LLM in process modeling is the study by Kourani, Berti, Schuster & van der Aalst (2024*a*) which has taken a look at quality and time, but not at costs; they also do not study the trade-off between relevant measures or identify a Pareto front. Therefore, we complemented this study with our additional perspective on modeling costs and there relationship. (5) Another important line of work is on data which can be used as benchmarks for modeling, e.g. Walter et al. (2014), Bellan et al. (2022), Sola et al. (2022). Not all of these datasets do contain a textual description of the domain to be modeled but could be added with generative synthetic approaches.

# 6 Conclusions and further research work

We presented the design for an experimental laboratory study of the process modeling abilities of LLM. Compared to a field study, always validity questions arise, e.g. regarding the representativeness of the standard modeling scenario. Although the used modeling scenarios are divers, we do not have a precise qualification for the validity of the samples. Since known evaluations only use more or less convenient samples, the results should be treated with caution.

In fact, it can be argued that the LLM and training datasets used are inherently problematic. Since many available foundation models are not open and transparent (Liesenfeld et al. (2023), Huang et al. (2024)), there is a chance that the foundational models used the known datasets for their own training. Formulated in more technical terms, if the well-known datasets are used for evaluation, the percentage of data leakage might be 100 percent. Without transparency, the comparison and evaluation of the abilities of LLM is almost uncontrollable and thus biased to an unknown extent. This problem is well-known in the context of the domain of protein prediction (Marx 2022). For an adequate evaluation of a system for protein prediction, biomolecules with unknown protein structures must be used. Otherwise, it cannot be excluded that the LLM is "cheating" because the correct solution is known and contained in the training data set. Hence, systematic studies of the generalization abilities of LLM are needed.

Furthermore, the introduced assumptions of the standard evaluation scenario can be relaxed. First, BPMN has a primary focus on the flow of activities and does not explicitly model data objects. However, it has been well known for decades that data is of major importance to understand business processes (Scheer 1994, Weske 2019). Hence, other modeling languages might be used to evaluate the results of an LLM.

Secondly, instead of a textual description, other modalities for providing information about the domain of modeling can be provided, e.g. images, video, or speech (Chvirova et al. 2024). In this direction, it might be interesting to set up a virtual laboratory environment to detect and model business processes, e.g. Knoch et al. (2020) already set up such scenarios, but do not use them for the evaluation of LLM modeling abilities.

Third, it might be interesting not only to use the abilities of an LLM for solving the standard modeling scenario. Instead process modeling may be assisted by different LLM, e.g. evaluating the solution, providing feedback, gaining data, improving the solution, translating natural language used for domain description, or interactively training the LLM. It is obvious that such an interactive modeling scenario ("agentic process modeling") is attractive but cannot easily be controlled in real-world or laboratory scenarios.

Fourth, the previous aspect already focuses on particular application domains of process modeling. Since the vocabulary and the workflows in process modeling are strongly dependent on the application domain, it might be necessary to focus on particular domains for evaluation, e.g. manufacturing, health, insurance. If a particular domain is focused, it might be necessary to fine-tune the LLM for this particular purpose.

In summary, a valid evaluation of the abilities of an LLM for process modeling is not only of practical relevance but naturally leads to several theoretically interesting questions for the field of process modeling.

# References

Beck, P.-D., Mahlmeister, T., Ifland, M. & Puppe, F. (2015), COCLAC - Feedback generation for combined UML class and activity diagram modeling tasks, *in* 'Proc. of 2nd Workshop "Automatische Bewertung von Programmieraufgaben" (ABP 2015)', Vol. 1496 of *CEUR-WS*.

Bellan, P., van der Aa, H., Dragoni, M., Ghidini, C. & Ponzetto, S. P. (2022), PET: an annotated dataset for process extraction from natural language text tasks, *in* C. Cabanillas, N. F. Garmann-Johnsen & A. Koschmider, eds, 'Business Process Management Workshops - BPM 2022 International Workshops, Münster, Germany, September 11-16, 2022, Revised Selected Papers', Vol. 460 of *Lecture Notes in Business Information Processing*, Springer, pp. 315–321.

Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A., Creel, K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P. W., Krass, M., Krishna, R., Kuditipudi, R., Kumar, A., Ladhak, F., Lee, M., Lee, T., Leskovec, J., Levent, I., Li, X. L., Li, X., Ma, T., Malik, A., Manning, C. D., Mirchandani, S., Mitchell, E., Munyikwa, Z., Nair, S., Narayan, A., Narayanan, D., Newman, B., Nie, A., Niebles, J. C., Nilforoshan, H., Nyarko, J., Ogut, G., Orr, L., Papadimitriou, I., Park, J. S., Piech, C., Portelance, E., Potts, C., Raghunathan, A., Reich, R., Ren, H., Rong, F., Roohani, Y., Ruiz, C., Ryan, J., Ré, C., Sadigh, D., Sagawa, S., Santhanam, K., Shih, A., Srinivasan, K., Tamkin, A., Taori, R., Thomas, A. W., Tramèr, F., Wang, R. E., Wang, W., Wu, B., Wu, J., Wu, Y., Xie, S. M., Yasunaga, M., You, J., Zaharia, M., Zhang, M., Zhang, T., Zhang, X., Zhang, Y., Zheng, L., Zhou, K. & Liang, P. (2022), 'On the opportunities and risks of foundation models', *arXiv preprint arXiv:2108.07258* .

Bunke, H. & Shearer, K. (1998), 'A graph distance metric based on the maximal common subgraph', *Pattern recognition letters* **19**(3-4), 255–259.

Chvirova, D., Egger, A., Fehrer, T., Kratsch, W., Röglinger, M., Wittmann, J. & Wördehoff, N. (2024), 'A multimedia dataset for object-centric business process mining in it asset management', *Data in Brief* **55**, 110716.

Dijkman, R., Dumas, M., van Dongen, B., Käärik, R. & Mendling, J. (2011), 'Similarity of business process models: Metrics and evaluation', *Information Systems* **36**(2), 498–516. Special Issue: Semantic Integration of Data, Multimedia, and Services.

Erdogan, L. E., Lee, N., Kim, S., Moon, S., Furuta, H., Anumanchipalli, G., Keutzer, K. & Gholami, A. (2025), 'Plan-and-act: Improving planning of agents for long-horizon tasks'.

Fill, H.-G., Fettke, P. & Köpke, J. (2023), 'Conceptual modeling and large language models: Impressions from first experiments with chatgpt', *Enterprise Modelling and*

*Information Systems Architectures (EMISAJ) – International Journal of Conceptual Modeling* **18**, 1–15.

Grandi, D., Jain, Y. P., Groom, A., Cramer, B. & McComb, C. (2024), 'Evaluating large language models for material selection', *Journal of Computing and Information Science in Engineering* **25**(2), 021004.

Gutschmidt, A. & Nast, B. (2025), Assessing model quality using large language models, *in* E. Paja, J. Zdravkovic, E. Kavakli & J. Stirna, eds, 'The Practice of Enterprise Modeling', Springer Nature Switzerland, Cham, pp. 105–122.

Heggset, M., Krogstie, J. & Wesenberg, H. (2015), Understanding model quality concerns when using process models in an industrial company, *in* K. Gaaloul, R. Schmidt, S. Nurcan, S. Guerreiro & Q. Ma, eds, 'Enterprise, Business-Process and Information Systems Modeling', Springer International Publishing, Cham, pp. 395–409.

Houy, C., Fettke, P. & Loos, P. (2012), Understanding understandability of conceptual models – what are we actually talking about?, *in* P. Atzeni, D. Cheung & S. Ram, eds, 'Conceptual Modeling', Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 64–77.

Hu, T. & Zhou, X.-H. (2024), 'Unveiling llm evaluation focused on metrics: Challenges and solutions'.

Huang, Y., Sun, L., Wang, H., Wu, S., Zhang, Q., Li, Y., Gao, C., Huang, Y., Lyu, W., Zhang, Y., Li, X., Liu, Z., Liu, Y., Wang, Y., Zhang, Z., Vidgen, B., Kailkhura, B., Xiong, C., Xiao, C., Li, C., Xing, E., Huang, F., Liu, H., Ji, H., Wang, H., Zhang, H., Yao, H., Kellis, M., Zitnik, M., Jiang, M., Bansal, M., Zou, J., Pei, J., Liu, J., Gao, J., Han, J., Zhao, J., Tang, J., Wang, J., Vanschoren, J., Mitchell, J., Shu, K., Xu, K., Chang, K.-W., He, L., Huang, L., Backes, M., Gong, N. Z., Yu, P. S., Chen, P.-Y., Gu, Q., Xu, R., Ying, R., Ji, S., Jana, S., Chen, T., Liu, T., Zhou, T., Wang, W., Li, X., Zhang, X., Wang, X., Xie, X., Chen, X., Wang, X., Liu, Y., Ye, Y., Cao, Y., Chen, Y. & Zhao, Y. (2024), 'Trustllm: Trustworthiness in large language models', *arXiv preprint arXiv:2401.05561* .

Joel, S., Wu, J. J. & Fard, F. H. (2024), 'A survey on llm-based code generation for low-resource and domain-specific programming languages'.

Kapoor, S., Stroebl, B., Siegel, Z. S., Nadgir, N. & Narayanan, A. (2024), 'Ai agents that matter', *arXiv preprint arXiv:2401.00001* .

Katz, M., Kokel, H., Srinivas, K. & Sohrabi Araghi, S. (2024), 'Thought of search: Planning with language models through the lens of efficiency', *Advances in Neural Information Processing Systems* **37**, 138491–138568.

Knoch, S., Ponpathirkoottam, S. & Schwartz, T. (2020), Video-to-model: Unsupervised trace extraction from videos for process discovery and conformance checking in manual assembly, *in* D. Fahland, C. Ghidini, J. Becker & M. Dumas, eds, 'Business Process Management - 18th International Conference, BPM 2020, Seville, Spain, September 13-18, 2020, Proceedings', Vol. 12168 of *Lecture Notes in Computer Science*, Springer, pp. 291–308.

Kourani, H., Berti, A., Hennrich, J., Kratsch, W., Weidlich, R., Li, C.-Y., Arslan, A., Schuster, D. & van der Aalst, W. M. P. (2024), 'Leveraging large language models for enhanced process model comprehension', *arXiv preprint arXiv:2408.08892* .

Kourani, H., Berti, A., Schuster, D. & van der Aalst, W. M. P. (2024*a*), 'Evaluating large language models on business process modeling: Framework, benchmark, and self-improvement analysis', *arXiv preprint arXiv:2402.00001* .

Kourani, H., Berti, A., Schuster, D. & van der Aalst, W. M. P. (2024*b*), Process modeling with large language models, *in* H. van der Aa, D. Bork, R. Schmidt & A. Sturm, eds, 'Enterprise, Business-Process and Information Systems Modeling. LNBIP 511', Springer Nature Switzerland, Cham, pp. 229–244.

Krogstie, J. (2016), *Quality in Business Process Modeling*, Springer.

Liesenfeld, A., Lopez, A. & Dingemanse, M. (2023), Opening up chatgpt: Tracking openness, transparency, and accountability in instruction-tuned text generators, *in* 'Proceedings of the 5th International Conference on Conversational User Interfaces', CUI '23, ACM, p. 1–6.

Marx, V. (2022), 'Method of the year: protein structure prediction', *Nature Methods* **19**(2), 5–10.

Mizrahi, M., Kaplan, G., Malkin, D., Dror, R., Shahaf, D. & Stanovsky, G. (2024), 'State of what art? a call for multi-prompt llm evaluation', *Transactions of the Association for Computational Linguistics* **12**, 933–949.

Norouzifar, A., Kourani, H., Dees, M. & van der Aalst, W. M. (2024), 'Bridging domain knowledge and process discovery using large language models', *arXiv preprint arXiv:2408.17316* .

Pavlicek, J., Pavlickova, P. & Naplava, P. (2019), Measures of quality in business process modeling, *in* R. Pergl, E. Babkin, R. Lock, P. Malyzhenkov & V. Merunka, eds, 'Enterprise and Organizational Modeling and Simulation', Springer International Publishing, Cham, pp. 146–155.

Sanchez-Ferreres, J., Delicado, L., Andaloussi, A. A., Burattin, A., Calderon-Ruiz, G., Weber, B., Carmona, J. & Padro, L. (2020), 'Supporting the Process of Learning and Teaching Process Models', *IEEE Transactions on Learning Technologies* **13**(3), 552–566.

Sarkar, A., Gordon, A. D., Negreanu, C., Poelitz, C., Ragavan, S. S. & Zorn, B. (2022), 'What is it like to program with artificial intelligence?'.

Scheer, A. (1994), 'ARIS toolset: A software product is born', *Inf. Syst.* **19**(8), 607–624.

Schneider, J., Meske, C. & Kuss, P. (2024), 'Foundation models', *Business & Information Systems Engineering* **66**(2), 221–231.

Schramm, J., Strickroth, S., Le, N.-T. & Pinkwart, N. (2012), Teaching UML skills to novice programmers using a sample solution based intelligent tutoring system, *in* 'Proc. of 25th International Florida Artificial Intelligence Research Society Conference (FLAIRS-25)', pp. 472–477.

Sola, D., Warmuth, C., Schäfer, B., Badakhshan, P., Rehse, J. & Kampik, T. (2022), SAP signavio academic models: A large process model dataset, *in* M. Montali, A. Senderovich & M. Weidlich, eds, 'Process Mining Workshops - ICPM 2022 International Workshops, Bozen-Bolzano, Italy, October 23-28, 2022, Revised Selected Papers', Vol. 468 of *Lecture Notes in Business Information Processing*, Springer, pp. 453–465.

Stein, K., Fišer, D., Hoffmann, J. & Koller, A. (2025), 'Automating the generation of prompts for llm-based action choice in pddl planning'.

Striewe, M. & Goedicke, M. (2014), Automated assessment of UML activity diagrams, *in* 'Proc. of 19th Annual Innovation and Technology in Computer Science Education Conference (ITiCSE 2014)', ACM, p. 336.

Thaler, T., Hake, P., Fettke, P. & Loos, P. (2014), Evaluating the evaluation of process matching techniques, *in* D. Kundisch, L. Suhl & L. Beckmann, eds, 'Multikonferenz Wirtschaftsinformatik, MKWI 2014, Paderborn, Germany, February 26-28, 2014', University of Paderborn, pp. 1600–1612.

Thaler, T., Houy, C., Fettke, P. & Loos, P. (2016), Automated assessment of process modeling exams: Basic ideas and prototypical implementation, *in* S. Betz & U. Reimer, eds, 'Modellierung 2016 - Workshopband. Workshop zur Modellierung in der Hochschullehre (MoHoL-2016)', Vol. 255 of *Lecture Notes in Informatics (LNI)*, Gesellschaft für Informatik (GI), Bonn, pp. 63–70.

Ullrich, M., Houy, C., Stottrop, T., Striewe, M., Willems, B., Fettke, P., Loos, P. & Oberweis, A. (2023), 'Automated assessment of conceptual models in education - a systematic literature review', *Enterprise Modelling and Information Systems Architectures - International Journal of Conceptual Modeling (EMISAJ)* **18**(2), 1–36.

Vidgof, M., Bachhofner, S. & Mendling, J. (2023), Large language models for business process management: Opportunities and challenges, *in* C. D. Francescomarino, A. Burattin, C. Janiesch & S. Sadiq, eds, 'Business Process Management Forum - BPM 2023 Forum, Utrecht, The Netherlands, September 11-15, 2023, Proceedings', Vol. 490 of *Lecture Notes in Business Information Processing*, Springer, pp. 107–123.

Voelter, M., Hadian, R., Kampik, T., Breitmayer, M. & Reichert, M. (2024), 'Leveraging generative ai for extracting process models from multimodal documents'.

von Rosing, M., White, S., Cummins, F. & de Man, H. (2015), Business process model and notation—bpmn, *in* M. von Rosing, A.-W. Scheer & H. von Scheel, eds, 'The Complete Business Process Handbook', Morgan Kaufmann, Boston, pp. 433–457.

Walter, J., Thaler, T., Ardalani, P., Fettke, P. & Loos, P. (2014), Development and usage of a process model corpus, *in* B. Thalheim, H. Jaakkola, Y. Kiyoki & N. Yoshida, eds, 'Information Modelling and Knowledge Bases XXVI, 24th International Conference on Information Modelling and Knowledge Bases (EJC 2014), Kiel, Germany, June 3-6, 2014', Vol. 272 of *Frontiers in Artificial Intelligence and Applications*, IOS Press, pp. 437–448.

Wang, L., Yang, F., Zhang, C., Lu, J., Qian, J., He, S., Zhao, P., Qiao, B., Huang, R., Qin, S., Su, Q., Ye, J., Zhang, Y., Lou, J.-G., Lin, Q., Rajmohan, S., Zhang, D. & Zhang, Q. (2025), 'Large action models: From inception to implementation'.

Weske, M. (2019), *Business Process Management - Concepts, Languages, Architectures, Third Edition*, Springer.

Westergaard, M., Fahland, D. & Stahl, C. (2013), Grade/CPN: A tool and temporal logic for testing colored Petri net models in teaching, *in* 'Proc. of 33rd International Conference on Application and Theory of Petri Nets and Other Models of Concurrency (Petri Nets 2012)', Vol. 8100, pp. 180–202.

Wiśniewski, P., Kluza, K. & Ligęza, A. (2019), Towards automated process modeling based on bpmn diagram composition, *in* F. Daniel, Q. Z. Sheng & H. Motahari, eds, 'Business Process Management Workshops', Springer International Publishing, Cham, pp. 507–513.

Zhang, J., Lan, T., Zhu, M., Liu, Z., Hoang, T., Kokane, S., Yao, W., Tan, J., Prabhakar, A., Chen, H., Liu, Z., Feng, Y., Awalgaonkar, T., Murthy, R., Hu, E., Chen, Z., Xu, R., Niebles, J. C., Heinecke, S., Wang, H., Savarese, S. & Xiong, C. (2024), 'xlam: A family of large action models to empower ai agent systems'.