
The Need for Tabular Representation Learning: An Industry Perspective

Alexandra Savelieva, Andreas Mueller, Avriella Floratou, Carlo Curino, Hiren Patel,
Jordan Henkel, Joyce Cahoon, Markus Weimer, Nellie Gustafsson, Richard Wydrowski,
Roman Batoukov, Shaleen Deep, Venkatesh Emani
first_name.last_name@microsoft.com

Abstract

1 The total addressable market for data applications has been estimated at \$70B. This
2 includes the \$11B market for data integration, which is estimated to grow at 25%
3 in the coming year; \$35B market for analytics, growing at 11%; and \$19B market
4 for business intelligence, growing at 8% [1]. Given this data-driven future and the
5 scale at which Microsoft operates, we survey PMs, engineers and researchers and
6 synthesize their opinions around extracting insights from tabular data at-scale. We
7 see three main areas where tabular representation learning (TRL) can be leveraged:

8 **Data insights.** Enabling real-time analytics is one of the key priorities for Mi-
9 crosoft's new intelligence platform [6] now that a converged environment exists
10 to house any type of data. TRL models can help expose column and table-level
11 semantic annotations, relationships between columns and between tables, and
12 advanced data patterns such as semantic-aware denial constraints [5].

13 **Data management.** From our internal workload telemetry, we know that 17.8% of
14 tabular data across our virtual clusters remain unaccessed [9]. From an external
15 perspective, leveraging telemetry from Azure Observability Platform, we observed
16 that out of the 10B+ metrics generated, less than 0.1% is used [3]. Bringing this
17 data to light requires sophisticated data discovery, data understanding and data
18 integration capabilities. We believe TRL models can play an important role on tasks
19 such as entity detection and deduplication, schema mapping, and data imputation.

20 **Data movement.** It is well-known that data movement remains a key bottleneck
21 in analytics [10]. In order to ensure that our users receive the best performance
22 possible, investments have been made in smart caching policies, like those involving
23 materialized views [4], as well as predicate operator pushdown [2]. Recent work
24 [8] predicts various structural and performance properties of queries by pre-training
25 encoder models with database workloads; but, the application of these strategies
26 fail to consider the underlying tabular data. With TRL models, we can jointly
27 pre-train tables with their query plans to enhance our understanding and ability to
28 characterize workloads, and thus further efforts in reducing data movement.

29 **Challenges and opportunities.** Existing tabular models are mostly trained on
30 Wikipedia tables and/or spreadsheets. However, in an enterprise setting, both the
31 customer data and their associated schema is often industry specific. Access to
32 the customer data is typically not possible due to privacy regulations [11, 7], thus
33 training TRL models on such data is often not possible. It remains an open question,
34 whether the existing TRL models can be successfully used on domain-specific
35 data. Along the same lines, the existence of large language models (LLMs) and
36 Microsoft's exclusive license to them, allows rapid prototyping of many applica-
37 tions even on top of tabular data. We encourage the community to provide studies
38 that compare the performance of LLMs and TRL on some of the tasks mentioned
39 above. Such systematic studies will be useful to application developers and product
40 teams that are looking to incorporate more ML-based capabilities.

41 **References**

- 42 [1] Analytics and business intelligence platforms. [https://www.gartner.com/reviews/](https://www.gartner.com/reviews/market/analytics-business-intelligence-platforms)
43 [market/analytics-business-intelligence-platforms](https://www.gartner.com/reviews/market/analytics-business-intelligence-platforms). Accessed: 2022-09-15.
- 44 [2] Using S3 select pushdown with Presto to improve performance. [https://](https://docs.aws.amazon.com/emr/latest/ReleaseGuide/emr-presto-s3select.html)
45 docs.aws.amazon.com/emr/latest/ReleaseGuide/emr-presto-s3select.html.
46 Accessed: 2022-11-06.
- 47 [3] Souren Aghajanyan, Roman Batoukov, and Jian Zhang. Signal Fabric—An AI-assisted
48 platform for knowledge discovery in dynamic system. Santa Clara, CA, May 2019. USENIX
49 Association.
- 50 [4] Rana Alotaibi, Bogdan Cautis, Alin Deutsch, Moustafa Latrache, Ioana Manolescu, and Yifei
51 Yang. Estocada: Towards scalable polystore systems. *Proceedings of the VLDB Endowment*,
52 13(12):2949–2952, 2020.
- 53 [5] Xu Chu, Ihab F Ilyas, and Paolo Papotti. Discovering denial constraints. *Proceedings of the*
54 *VLDB Endowment*, 6(13):1498–1509, 2013.
- 55 [6] Rohan Kumar. Introducing the Microsoft Intelligent Data Platform.
56 [https://azure.microsoft.com/en-us/blog/](https://azure.microsoft.com/en-us/blog/introducing-the-microsoft-intelligent-data-platform/)
57 [introducing-the-microsoft-intelligent-data-platform/](https://azure.microsoft.com/en-us/blog/introducing-the-microsoft-intelligent-data-platform/). Accessed: 2022-09-15.
- 58 [7] State of California Dept of Justice. California consumer privacy act.
59 <https://oag.ca.gov/privacy/ccpa>. Accessed: 2022-11-10.
- 60 [8] Debjyoti Paul, Jie Cao, Feifei Li, and Vivek Srikumar. Database workload characterization with
61 query plan encoders. *Proceedings of the VLDB Endowment*, 15(4):923–935, 2021.
- 62 [9] Conor Power, Hiren Patel, Alekh Jindal, Jyoti Leeka, Bob Jenkins, Michael Rys, Ed Triou,
63 Dexin Zhu, Lucky Katahanas, Chakrapani Bhat Talapady, et al. The Cosmos big data platform
64 at Microsoft: Over a decade of progress and a decade to look forward. *Proceedings of the*
65 *VLDB Endowment*, 14(12):3148–3161, 2021.
- 66 [10] Karla Saur, Tara Mirmira, Konstantinos Karanasos, and Jesús Camacho-Rodríguez.
67 Containerized execution of UDFs: An experimental evaluation. *Proceedings of the VLDB*
68 *Endowment*, 15(12):3158–3171, 2022.
- 69 [11] European Union. General data protection regulation.
70 <https://eur-lex.europa.eu/legal-content/EN/ALL/?uri=celex>. Accessed:
71 2022-11-10.