

# Michael Diskin

✉ [michael.s.diskin@gmail.com](mailto:michael.s.diskin@gmail.com) 📍 Europe 🗣️ [yhn112](#) 🌐 [yhn112](#)

## Experience

---

- |   |   |
|---|---|
| <b>Wildberries</b> , Head of LLM R&D <ul style="list-style-type: none"><li>Built LLM R&amp;D org from 0 to 30+ engineers (4–5 teams); manage a 500+ GPU cluster and define platform strategy for LLM and embeddings company-wide.</li><li>Shipped universal text embeddings (retrieval, ranking, classification) and market-place-scale MT translating tens of millions of product listings into 10+ languages, including low-resource.</li><li>Deployed optimized LLM serving with use-case routing, batching, and quantization — cut GPU costs by 40%+ while meeting latency SLAs under peak traffic.</li><li>Built RAG-powered assistants (internal and seller-facing) with safety guardrails, evaluation framework, and gold-standard regression datasets.</li><li>Established research-to-production operating model: experiment standards, evaluation gates, reliable releases — reduced iteration cycles from weeks to days.</li></ul> | Moscow, Russia<br>2024 – present<br>2 years |
| <b>Brask AI</b> , Senior Research Engineer <ul style="list-style-type: none"><li>Led R&amp;D on a lip-sync model for AI video dubbing: research, model iteration, production integration.</li><li>Improved robustness across diverse speakers and conditions; defined quality metrics and failure-analysis pipeline with product and engineering.</li></ul>   | Tbilisi, Georgia<br>2022 – 2023<br>1 year   |
| <b>Yandex Research</b> , Research Scientist <ul style="list-style-type: none"><li>Research on efficient and distributed training for large models; 5 papers at NeurIPS, ICML, and ICLR.</li><li>Led pre-training of a fully open-source Russian BERT-class language model; co-authored the Hivemind decentralized training library.</li><li>Created an evaluation benchmark for graph neural networks under heterophily (460+ citations, ICLR 2023).</li></ul>  | Moscow, Russia<br>2021 – 2022<br>1 year     |
| <b>Huawei</b> , Research Engineer <ul style="list-style-type: none"><li>Computer vision research: optical flow and depth estimation for image-based localization.</li></ul>   | Moscow, Russia<br>2020 – 2021<br>1 year     |
| <b>Early-stage startups</b> , ML / Software Engineer <ul style="list-style-type: none"><li>Built ML-powered backend services and data pipelines from scratch in small teams; end-to-end ownership from prototyping to deployment.</li></ul>   | Moscow, Russia<br>2019 – 2020<br>1 year     |
| <b>Yandex</b> , Software Engineering Intern <ul style="list-style-type: none"><li>Large-scale analytics over multi-terabyte log data using internal MapReduce infrastructure (YT).</li></ul>  | Moscow, Russia<br>2017 – 2018<br>1 year     |

## Education

---

- |   |                               |
|---|-------------------------------|
| <b>HSE University</b> , Computer Science                                    | Moscow, Russia<br>2022 – 2024 |
| <b>Yandex School of Data Analysis</b> , Machine Learning (Graduate program) | 2019 – 2021                   |
| <b>HSE University</b> , Computer Science                                    | Moscow, Russia<br>2014 – 2019 |

## Skills

---

**ML & AI**

**Frameworks**

**Data & Eval**

**Infra & MLOps**

**Languages**

## Awards

---

<b>HSE FCS Scholarship for Research Excellence</b> HSE University	2024
<b>Tinkoff Education Scholarship</b> Tinkoff	2023
<b>Xeek.ai "Put it on the Map!" — 2nd place (100+ teams)</b> Cash prize Xeek.ai	2020
<b>Kaggle "Recursion Cellular Image Classification" — 13th of 800+ teams</b> Kaggle	2019
<b>Huawei Image Inpainting Hackathon — 2nd place (100+ teams)</b> Cash prize Huawei	2019
<b>International Data Analysis Olympiad — 16th of 1000+ participants</b> IDAO	2019

## Languages

---

**Russian**

Native speaker

**English**

Fluent