

# Ian Yu

Pragmatic ML Engineer in the Age of Generative AI | 647 972 4689 | ian.yu@arc.com.co | LinkedIn | GitHub

## Highlights

- Machine Learning Engineer with 5+ years of experience designing and deploying production-grade ML systems in the eCommerce industry. Deployed 5 agentic projects to production and contributed to 2 open research initiatives on LLMs.
- Proficient in production-grade NLP, agentic workflows, RAG, tool-use, MCPs, OCR, MLOps, analytics data engineering, excel at designing-prototype to implement-at-scale transition
- Experienced in consultancy, communication with various levels of stakeholders, product thinking and strategize with users, specialized in utilizing subject matter knowledge with AI

## Relevant Experience

<b>LLM Engineer</b> , Pixomondo (Sony Subsidiary) – Toronto, ON	Nov 2025 – Present
<ul style="list-style-type: none"> <li>Assigned to a high impact innovation platform project with most artifacts built from scratch</li> <li>Built custom system discovery and registry flow for internal tools on an event-driven platform, enabling agent tool use over 20+ system services</li> <li>Designing and implementing event-driven internal tool use evaluation, modification suggestion, and continual learning</li> </ul>	
<b>Machine Learning Engineer</b> , Aggregate Intellect, Contract – Toronto, ON	Apr 2024 – Present
<ul style="list-style-type: none"> <li>For a RAG ChatBot client with high-fidelity requirements in legal fields, redesigned agentic workflows with 10x lower latency, 3-5x lower cost, while maintaining same output quality</li> <li>Architect and implemented agentic systems for a legal tech company, including document intelligence (Wordx and PDFs), RAG conversations, and memory.</li> </ul>	
<b>Machine Learning Engineer</b> , RezolveAI (formerly Groupby Inc) – Toronto, ON	Aug 2022 – Oct 2025
<ul style="list-style-type: none"> <li>Led, designed, and productionized Enrich AI, a major product with agentic workflow that manages product taxonomy, data strategy, and information enrichment, cutting full-service project timelines from weeks to hours.</li> <li>Built autotuning text clustering system with custom loss and outlier handling for efficient AI reasoning reviews, deployed via Ray and Argo on GKE. Lowered orders of magnitude compute while increasing quality by 70%</li> <li>Revamped AI in a patent-pending automotive fitment solution using RAG, achieving &gt;1000 QPS and maintaining &lt;1s latency at 99th percentile.</li> </ul>	
<b>Data Strategy Analyst</b> , Groupby Inc – Toronto, ON	Apr 2021 – Jul 2022
<ul style="list-style-type: none"> <li>Created internal packages for client-specific data enrichment analysis, boosting client satisfaction by 50% and reducing revisions by 40%. Automated NLP/ML processes enhancing data observability and collaboration, improving team productivity by 50%.</li> <li>Introduced weakly-supervised text classification system to withhold null tasks from data annotators, saved \$100,000 USD throughout its lifetime.</li> </ul>	

## Open Contribution

<b>The BigScience ROOTS Corpus: A 1.6TB Composite Multilingual Dataset</b> [arxiv]	2021 – 2022
<ul style="list-style-type: none"> <li>Led a team of 10 individual contributors to minimize Personally Identifiable Information leakage. Analyzed and defined data filtering parameters to ensure corpus quality, safeness, and ethics</li> <li>NeurIPS 2022, Datasets and Benchmarks Track</li> </ul>	
<b>SantaCoder: don't reach for the stars!</b> [arxiv]	2022 – 2022
<ul style="list-style-type: none"> <li>Advised on data annotation best practices and framework to facilitate PII detection effort for the Stack corpus</li> </ul>	

## Technologies

<b>Machine Learning:</b> PyTorch, Transformers, Huggingface, Scikit-learn, SciPy, SpaCy, Snorkel, Vector Databases, Arrows, Polars, SageMaker, Vertex AI, Ray, Kubeflow, MLflow, DSPy
<b>Engineering:</b> Python, JavaScript, TypeScript, SQL, NoSQL, Bash, Linux, Docker, Kubernetes, Argo, CI/CD, gRPC, REST, GraphQL, OpenAPI, Locust, AWS, GCP