

Subject: Create your own AI Engineering Team
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Saturday, February 17, 2024 5:31:46 PM GMT+05:30
Date Received: Saturday, February 17, 2024 5:31:47 PM GMT+05:30

February 17, 2024 | [Read Online](#)

Create your own AI Engineering Team

learn all about roles and skills to boost your business



Welcome to the 2,138 new members this week! **nocode.ai** now has 38,728 subscribers

Having a squad of AI Engineers will make a huge difference for your organization. They will be able to create custom AI solutions leveraging LLMs and multiple components essential for creating Generative AI apps. All solutions custom for your own business.

Today I'll cover:

- The Power of Custom AI Solutions
- The roles needed
- How to attract talent
- Essential skills
- Team structure
- Talent retention strategies

Let's Dive In!

The Power of Custom AI Solutions

Imagine a suite of AI-driven solutions tailor-made for your business. With a proficient team of AI engineers, the possibilities extend far beyond conventional pre-built solutions. Those solutions offer convenience but they often lack the nuance and personalization needed to truly address your unique business challenges.

To truly embrace AI, you need to embrace in-house AI development.

A dedicated team can craft custom solutions tailored to your specific workflows, data, and objectives. Imagine:

- **Intuitive, Chat-Based Document Search:** Say goodbye to searching through endless documents. An AI-powered search tool, accessible through a simple chat interface, can answer your questions and pull relevant information instantly.
- **Autonomous Task Management:** Repetitive tasks bog down your team. Build AI agents that handle them seamlessly, freeing up human talent for more strategic work. Think scheduling, data entry, or even customer service interactions.
- **Personalized Customer Experiences:** Delight your customers with AI-driven chatbots that answer their questions, recommend products, and resolve issues efficiently. Personalization goes beyond simple interactions, with AI analyzing customer data to tailor experiences and increase satisfaction.
- **Content Creation on Autopilot:** Generate high-quality reports, marketing materials, or even social media posts using AI algorithms trained on your brand voice and style. Save time and resources while producing engaging content that resonates with your audience.

These are just a few examples. The possibilities are endless, limited only by your imagination and the expertise of your AI engineering team.

Building Your Dream Team: The Road to AI Success

Excited to tap into the power of custom AI? Let's navigate the journey to assembling your dream AI team, where expertise meets innovation.

Essential Roles

In the rapidly evolving AI landscape, the shift towards leveraging LLMs has transformed the composition and focus of AI teams:

Traditional ML Team Roles:

- **Machine Learning Engineers** build and optimize AI models.
- **Data Scientists** ensure model accuracy through data analysis.
- **Software Engineers** integrate AI models with existing systems.
- **UX/UI Designers** develop user-friendly interfaces for AI tools.

New Era Adjustments:

- **AI Engineers:** Now, this role merges the responsibilities of Machine Learning Engineers and Software Engineers. These professionals select, integrate, and potentially fine-tune pre-trained LLMs, focusing on application and customization rather than building models from scratch. Their expertise in LLMs simplifies connecting AI to systems and crafting effective prompts.
- **UX/UI Designers:** Their role evolves to address more complex human-AI interactions, ensuring that interfaces are not only intuitive but also capable of facilitating sophisticated engagements like conversational UIs and adaptive feedback mechanisms.

Mastering the Essential Skills for AI Engineers

In this issue, I dig into the fresh and exciting world of AI Engineering. I'll chat about what an AI Engineer really does, and I'll lay out what it takes for you to become one yourself.

www.nocode.ai/the-ai-engineer-role

The Impact of Pre-trained LLMs: Pre-trained LLMs have democratized AI by making powerful tools accessible without the need for extensive data or computational resources. This shift allows AI teams to concentrate on application, integration, and customization, speeding up development and innovation across sectors. It represents a significant step towards AI solutions that are more adaptable, intuitive, and personalized, underscoring a new era in AI technology's democratization.

Data Scientists are not becoming obsolete; their skills in interpreting AI data, fine-tuning models, and is still important. My recommendation though is to find AI Engineers who should also know the basics of model tuning to complement this expertise.

With LLMs, we have a leaner team and a more focused approach.

everyone should master ai

consumers of AI

- apply AI to business needs
- leverage pre-built services and ai solutions
- master prompting ai models

Examples: marketing, analyst, seller, support rep, product manager

creators of AI

- build end-to-end ai solutions
- work on ai technology directly
- Require ai expertise

Examples: ai engineers, developers, data scientists



consumers provide application insights to guide development



with no-code ai tools, **everyone can become an AI Creator**

the ai bootcamp

Everyone should learn AI

Attracting Top Talent

The competition for AI talent is fierce, but with the right strategy, you can attract the best. Here are some prime hunting grounds:

- **AI-focused job boards:** Platforms like LinkedIn, Hired, and Indeed feature dedicated sections for AI job listings.
- **Professional networking sites:** LinkedIn excels as a resource for finding AI professionals in industry-specific groups and communities.
- **Technical Community:** GitHub serves as a prime spot for identifying talent through AI projects and contributions.
- **University career fairs and programs:** Engage with top universities known for AI programs, which host career fairs and have job boards for graduates.
- **AI conferences and meetups:** These gatherings are ideal for networking with potential hires and promoting your company culture.

- **Employee referrals:** Utilize your team's network by incentivizing referrals to attract compatible and skilled candidates.

Salary Ranges: Competitive Offers for Top Talent

Compensation plays a crucial role in attracting and retaining AI talent. While specific figures can vary based on location, experience, and specialization, here's a general guide:

Level/Experience	United States	Europe	Asia
Entry-level (0-2 years)	\$90,000 - \$120,000	€50,000 - €70,000	\$30,000 - \$60,000
Mid-level (2-5 years)	\$120,000 - \$160,000	€70,000 - €100,000	\$60,000 - \$100,000
Senior-level (5+ years)	\$160,000 - \$210,000	€100,000 - €150,000	\$100,000 - \$150,000
Lead/Principal (8+ years)	\$210,000+	€150,000+	\$150,000+

Estimate of salaries for AI Engineers

Remember, competitive salary shouldn't be your only tool. Offer comprehensive benefits packages, opportunities for professional development, and a work environment that fosters creativity and innovation.

Building a Cohesive Unit

Harnessing AI's power starts with assembling a dream team that combines technical skills with a collaborative spirit. Here's how to build a cohesive unit:

- **Clear roles and responsibilities** avoid confusion and ensure everyone is working towards shared goals.
- **Open communication and collaboration** foster knowledge-sharing and problem-solving across disciplines, especially with other SMEs on the business side.
- **Continuous learning and development** keep your team at the forefront of the ever-evolving AI landscape.

Creating a balanced AI team structure is essential for organizations looking to leverage artificial intelligence effectively. Below is a recommended table

outlining the optimal number of AI Engineers relative to overall organizational size, and guidance on when to appoint a Head of AI.

Organization Size	Recommended Number of AI Engineers	When to Appoint a Head of AI
Small (1-50 employees)	1-2 AI Engineers	At 2 AI Engineers
Medium (51-200 employees)	3-5 AI Engineers	At 3 AI Engineers
Large (201-1000 employees)	6-10 AI Engineers	At 6 AI Engineers
Enterprise (1000+ employees)	10+ AI Engineers	At 10 AI Engineers

The optimal number of AI Engineers based on Organization Size

The importance of the Head of AI role

The **Head of AI** is responsible for the overall strategy, development, and implementation of AI initiatives within an organization or group.

The primary focus is to ensure that AI is being used effectively to drive business value and achieve the organization's goals. You need AI leaders who understand that AI projects require an Agile AI, quick iteration methodology where you start small and grow the implementations with validation over time.

Not like this...



...instead like this!



How to Develop AI Solutions

Responsibilities:

1. Lead ai teams
2. Driving AI research and development
3. ensuring ethical and responsible use of AI
4. communicating about ai



Armand Ruiz • You

Director AI at IBM

1w •

• • •

Every company or department of >30 people needs a Head of AI

To fully harness the power of AI, businesses need talent and expertise at every level:

- 🧠 Head of AI - An executive to craft your AI strategy and governance. They align AI with business goals.
- 👨‍💻 AI Engineer - Hands-on practitioners to operationalize AI systems end-to-end. They make it work for your needs.
- 🎓 Organization-wide training - Equip all knowledge workers with AI literacy and skills to apply AI tools daily. Democratization is key.
- 🤝 External AI partners - Complement internal capabilities by partnering with AI consultancies, agencies, and vendors.
- 🚀 C-suite sponsorship - Get buy-in from senior leadership to accelerate adoption. Give teams executive air cover.

You need a diverse blend of skill sets to transform.

AI is no longer optional - it's a competitive necessity.

I talk in more detail about the role of the Head of AI here:

<https://lnkd.in/dkv2U7AY>

The Importance of a Head of AI

Keeping Your AI Heroes Happy

Maximizing the potential of AI depends on the happiness and engagement of your team. Here's a streamlined strategy to ensure your AI professionals remain motivated and innovative:

- **Regular feedback and recognition** show appreciation for their contributions and foster motivation.
- **Opportunities for career growth and advancement** keep them engaged and challenged.
- **Work-life balance** enables them to perform at their best while maintaining personal well-being.

Category	Tactical Actions
Recognition	Monthly shout-outs for achievements, Performance-based bonuses or rewards
Professional Development	Sponsorship for AI conferences, Access to online courses and workshops
Career Advancement	Clear career paths with milestones, Regular career development meetings
Work-Life Balance	Flexible working hours, Remote work options
Community Engagement	Encourage participation in hackathons, Support contributions to open source projects

Tactical actions for enhancing your AI team's motivation

With the right team in place, your organization can unlock the immense potential of custom AI solutions. I hope you enjoyed today's edition.

Enjoy the weekend folks,

Armand

Whenever you're ready, there are FREE 2 ways to learn more about AI with me:

1. [The 15-day Generative AI course](#): Join my 15-day Generative AI email course, and learn with **just 5 minutes a day**. You'll receive concise daily lessons focused on practical business applications. It is perfect for quickly learning and applying core AI concepts. 10,000+ Business Professionals are already learning with it.

2. [The AI Bootcamp](#): For those looking to go deeper, join a full bootcamp with 50+ videos, 15 practice exercises, and a community to all learn together.



Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beejiiv

Subject: Creating an AI-Ready Culture
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Saturday, February 10, 2024 5:31:32 PM GMT+05:30
Date Received: Saturday, February 10, 2024 5:31:34 PM GMT+05:30

February 10, 2024 | [Read Online](#)

Creating an AI-Ready Culture

Adopt AI to truly Power your Business



Welcome to the 584 new members this week! **nocode.ai** now has 36,847 subscribers

The rise of AI isn't science fiction anymore. It's disrupting industries, automating tasks, and reshaping the future of work. But a harsh reality bites - most companies aren't ready for this revolution. They're like sleepwalkers stumbling towards a cliff, oblivious to the impending fall.

Don't be one of them.

Instead, wake up, and take proactive steps to prepare your business for the AI wave. By investing in three key areas - culture, change management, and rapid innovation - you can not only survive but thrive in this new era.

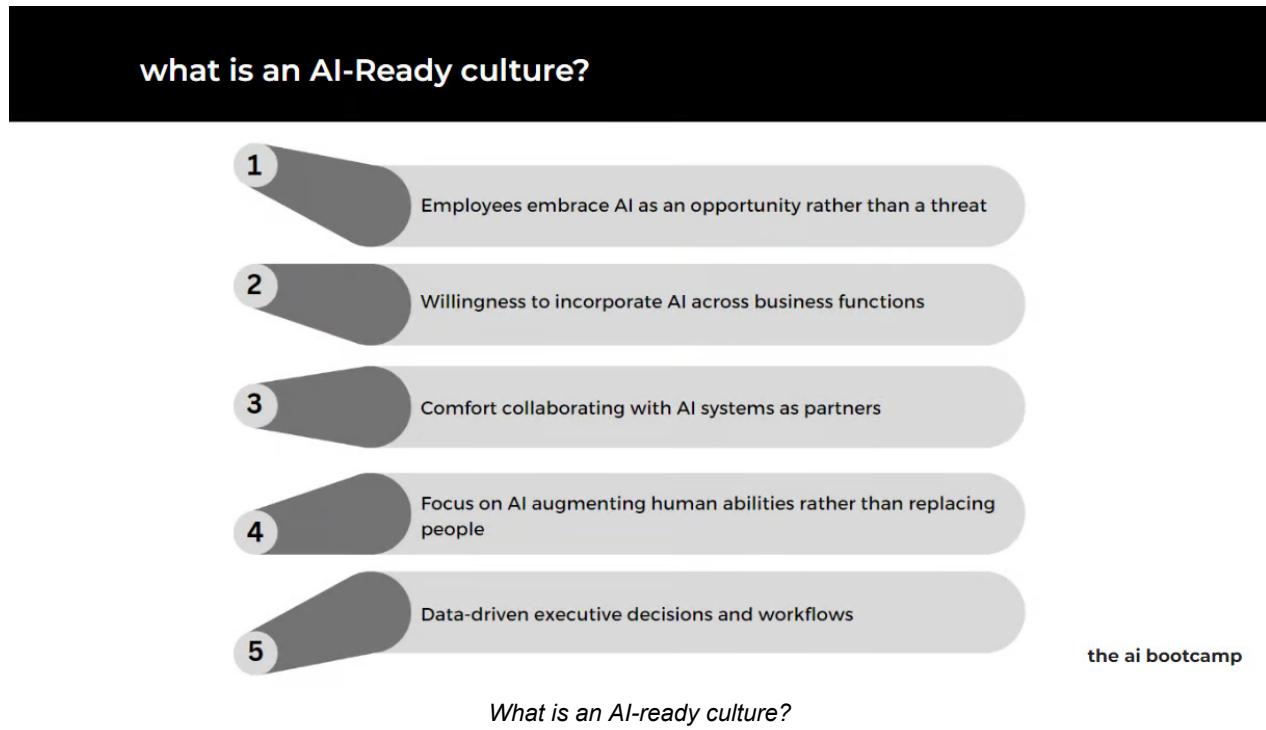
Today I'll cover:

- What is an AI-ready culture?
- Why companies Need to change their Culture to achieve success with AI
- How to build an AI-ready culture
- Benefits of an AI-ready culture

Let's Dive In!

What is an AI-ready culture?

AI isn't just a technology; it's a fundamental shift in mindset. Ditch the fear and embrace the potential. Foster a culture of **curiosity, learning, and adaptation**. Encourage employees to experiment with AI tools, attend workshops, and ask questions. Remember, a skilled workforce, not just AI itself, is your true competitive advantage.



Why companies Need to change their Culture to achieve success with AI

Don't wait for perfect solutions. Embrace agile methodologies and **experimentation** as your guiding principles. Create dedicated innovation labs, incentivize creative thinking, and reward calculated risks. Encourage collaboration between tech-savvy individuals and domain experts. Remember, **fast iteration, not stagnant perfection, is the key to AI success**.

Twelve months of sustained effort in these areas can make a world of difference. You'll foster a culture that welcomes AI, navigates change smoothly, and builds an environment where innovation thrives. The result? A future where your business isn't just surviving, but **leading the charge in the AI revolution**.

How to build an AI-ready culture and manage change

AI integration isn't a magic bullet. Expect bumps along the road. Proactively manage organizational change by being transparent, communicating effectively, and addressing

employee concerns. Train your workforce on new skills, invest in reskilling programs, and offer support throughout the transition. Remember, **happy and empowered employees are your biggest AI allies.**

building an ai-ready culture

Communicate an inspiring AI vision and strategy from leadership

Provide organization-wide AI literacy training and workshops

Incentivize employees to find opportunities to incorporate AI in their workflows

Develop dedicated roles like AI trainers, evangelists, and coaches

Celebrate successes and learn from failures through post-mortem reviews

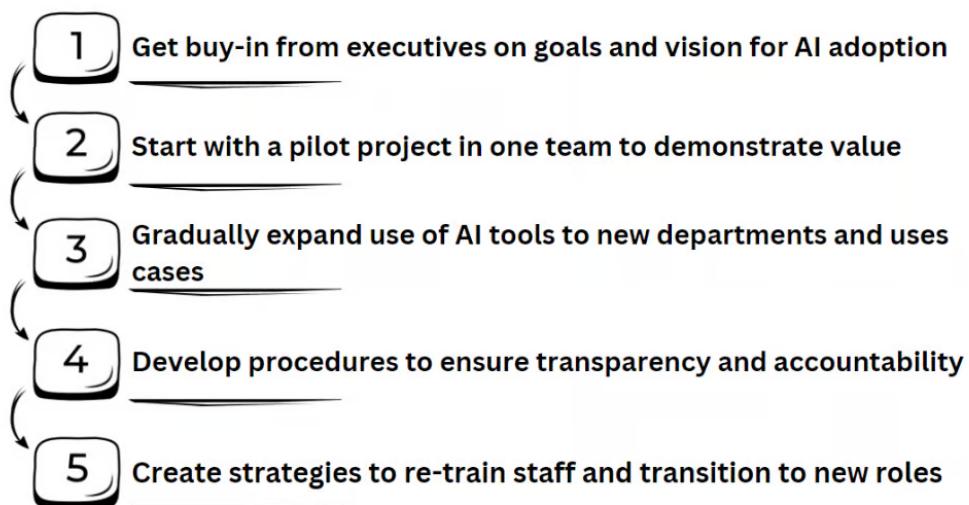
Survey employees regularly to gauge AI acceptance and continuously improve

the ai bootcamp

Building an AI-ready culture isn't an overnight sprint, it's a marathon. Here are some steps to get you started:

- **Leadership buy-in:** Secure support from top management to advocate for and champion AI initiatives.
- **Communication and education:** Educate employees about AI, its potential benefits, and how it will impact their roles. Address concerns openly and honestly.
- **Invest in training and development:** Equip employees with the skills needed to work with AI, including data literacy, critical thinking, and problem-solving.
- **Empowerment and experimentation:** Encourage employees to experiment with AI tools and suggest potential applications. Celebrate successes and learn from failures.
- **Create a feedback loop:** Gather feedback from employees on their experiences with AI and use it to continuously improve your approach.

Managing Organizational Change



the ai bootcamp

How to Manage Organizational Change

Benefits of an AI-ready culture

The benefits of fostering an AI-ready culture extend far beyond simply implementing AI tools. It can lead to:

- **Increased innovation:** A culture that encourages experimentation and embraces new ideas fosters a breeding ground for innovative AI solutions.
- **Improved decision-making:** Data-driven decisions based on AI insights can lead to better business outcomes.
- **Enhanced productivity:** AI can automate routine tasks, freeing up employees to focus on more strategic and creative work.
- **Stronger employee engagement:** When employees feel empowered and involved in AI initiatives, they're more likely to be engaged and motivated.
- **Attract and retain top talent:** A forward-thinking culture that embraces AI is attractive to the best and brightest minds.

benefits of an AI-Ready culture

1

enhanced workforce performance: AI adoption leads to increased employee productivity, satisfaction, and retention.

2

acceleration in innovation: Companies benefit from a faster pace of intelligent automation and AI solution implementation.

3

data-driven advantage: Adopting AI promotes better decision-making and gives a competitive edge over hesitant rivals.



the ai bootcamp

Still hesitant? Consider this:

- **Your competitors are likely making moves.** Don't let them gain an insurmountable advantage.
- **AI can unlock immense value.** From smarter operations to personalized customer experiences, the benefits are real.
- **Ignoring AI won't make it go away.** It's time to embrace the future, not run from it.

Remember, the choice is yours. Sleepwalk into oblivion or step boldly into the exciting, AI-powered future. The clock is ticking - will you answer the call?

Enjoy the weekend folks,

Armand

Whenever you're ready, there are FREE 2 ways to learn more about AI with me:

1. [**The 15-day Generative AI course**](#): Join my 15-day Generative AI email course, and learn with **just 5 minutes a day**. You'll receive concise daily lessons

focused on practical business applications. It is perfect for quickly learning and applying core AI concepts. 10,000+ Business Professionals are already learning with it.

2. **The AI Bootcamp**: For those looking to go deeper, join a full bootcamp with 50+ videos, 15 practice exercises, and a community to all learn together.



Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beejiiv

Subject: How to evaluate an LLM?
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Saturday, February 3, 2024 5:31:41 PM GMT+05:30
Date Received: Saturday, February 3, 2024 5:31:42 PM GMT+05:30

February 03, 2024 | [Read Online](#)

How to evaluate an LLM?

Metrics matter



Welcome to the 970 new members this week! **nocode.ai** now has 36,428 subscribers

Discover the critical aspects of evaluating LLMs to understand their capabilities and limitations. This guide outlines essential metrics and methodologies for a comprehensive assessment, ensuring the responsible development and deployment of Generative AI technology.

Today I'll cover:

- Why we need to evaluate LLMs
- Challenges in Evaluating LLMs
- Existing evaluation techniques
- Key Metrics for LLM Evaluation
- Future Directions in LLM Evaluation

Let's Dive In!

Why we need to evaluate LLMs

Evaluating LLMs is crucial as they increasingly underpin applications offering personalized recommendations, data translation, and summarization, among others.

With the proliferation of LLM applications, accurately measuring their performance becomes essential due to the limited availability of user feedback and the high cost and logistical challenges of human labeling.

Additionally, the complexity of LLM applications can make debugging difficult. Leveraging LLMs themselves for automated evaluation offers a promising solution to these challenges, ensuring that evaluations are both reliable and scalable.

Challenges in Evaluating LLMs

Evaluating LLMs involves addressing the subjective nature of language and the technical complexity of the models, alongside ensuring fairness and mitigating biases. As AI technology rapidly advances, evaluation methods must adapt to remain effective and ethical, demanding ongoing research and a balanced approach to meet these evolving challenges.

Here are 4 major hurdles we face:

1. **Biased Data, Biased Outputs:** Contaminated training data leads to unfair or inaccurate model responses. Identifying and fixing these biases in data and models is crucial.
2. **Beyond Fluency, Lies Understanding:** Metrics like perplexity focus on predicting the next word, not necessarily true comprehension. We need measures that capture deeper language understanding.
3. **Humans Can Be Flawed Evaluators:** Subjectivity and biases from human judges can skew evaluation results. Diverse evaluators, clear criteria, and proper training are essential.
4. **Real World Reality Check:** LLMs excel in controlled settings, but how do they perform in messy, real-world situations? Evaluation needs to reflect true-world complexities.

There is a survey that provides an in-depth look at evaluating LLMs, highlighting the importance of thorough assessments across knowledge, alignment, and safety to harness their benefits responsibly while mitigating risks. Here's the [link](#)

Evaluating Large Language Models: A Comprehensive Survey

Zishan Guo*, Renren Jin*, Chuang Liu*, Yufei Huang, Dan Shi, Supryadi Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, Deyi Xiong†

Tianjin University

{guozishan, rrjin, liuc_09, yuki_731, shidan, supryadi}@tju.edu.cn
{linhaoyu, yan_liu, jiaxuanlee, xbj1355, dyxiong}@tju.edu.cn

Abstract

Large language models (LLMs) have demonstrated remarkable capabilities across a broad spectrum of tasks. They have attracted significant attention and been deployed in numerous downstream applications. Nevertheless, akin to a double-edged sword, LLMs also present potential risks. They could suffer from private data leaks or yield inappropriate, harmful, or misleading content. Additionally, the rapid progress of LLMs raises concerns about the potential emergence of superintelligent systems without adequate safeguards.

Paper: Evaluating Large Language Models: A Comprehensive Survey

Existing evaluation techniques

Despite the challenges, researchers and developers have devised various techniques to assess LLMs' capabilities. Here are some prominent approaches:

- 1. Benchmark Datasets:** Standardized tasks like question answering (SQuAD), natural language inference (MNLI), and summarization (CNN/Daily Mail) offer controlled environments to compare LLM performance.
- 2. Automatic Metrics:** Metrics like BLEU score and ROUGE measure fluency and grammatical correctness by comparing the generated text to human-written references. However, they often prioritize surface-level similarity over deeper understanding.
- 3. Human Evaluation:** Crowdsourcing platforms and expert panels provide qualitative assessments of factors like coherence, creativity, and factual accuracy. This approach is subjective but offers valuable insights into aspects beyond numerical scores.
- 4. Adversarial Evaluation:** Crafting inputs designed to mislead LLMs helps expose vulnerabilities and evaluate their robustness against malicious attacks. This technique highlights potential safety concerns and guides LLM development.

5. Intrinsic Evaluation: Techniques like probing and introspection aim to assess an LLM's internal knowledge representations and reasoning processes. This offers glimpses into how the model "thinks" but is still under development.

These techniques provide valuable tools for LLM evaluation, but no single approach offers a complete picture. A multifaceted approach combining diverse techniques and addressing existing challenges will be crucial for building truly reliable and trustworthy LLMs.

The [HuggingFace Open LLM Leaderboard](#) aims to track, rank, and evaluate open LLMs and chatbots.

The screenshot shows the HuggingFace Open LLM Leaderboard interface. At the top, there is a header with the title "Open LLM Leaderboard" and a sub-header explaining the purpose: "The 🤗 Open LLM Leaderboard aims to track, rank and evaluate open LLMs and chatbots. Submit a model for automated evaluation on the 🤗 GPU cluster on the "Submit" page! The leaderboard's backend runs the great Eleuther AI Language Model Evaluation Harness - read more details in the "About" page!" Below the header are several filter and search options. On the left, there is a search bar and a section to "Select columns to show" with checkboxes for Average, ARC, HellaSwag, MMLU, TruthfulQA, Precision, Hub License, #Params (B), Hub, Model sha, and Type. There is also a checkbox for "Show gated/private/deleted models". On the right, there are sections for "Model types" (checkboxes for pretrained, fine-tuned, instruction-tuned, RL-tuned, and Unknown), "Precision" (checkboxes for torch.float16, torch.bfloat16, torch.float32, 8bit, 4bit, and GPTQ), and "Model sizes" (checkboxes for Unknown, <1.5B, ~3B, ~7B, ~13B, ~3SB, and 60B+). The main area is a table listing various LLM models with their metrics. The columns are T (Type), Model, Average, ARC, HellaSwag, MMLU, and TruthfulQA. The table includes rows for models like ICBU-NPU/FashionGPT-70B-V1.2, sequelbox/StellarBright, Riiid/sheep-duck-llama-2-70b-v1.1, AIDC-ai-business/Marcoroni-70B-v1, ICBU-NPU/FashionGPT-70B-V1.1, adonlee/LLaMA_2_70B_LoRA, uni-tianyan/Uni-TianYan, Riiid/sheep-duck-llama-2, Riiid/sheep-duck-llama-2, fangloveskari/ORCA_LLaMA_70B_QLoRA, ICBU-NPU/FashionGPT-70B-V3, and oh-yeontaek/llama-2-70B-LoRA-assemble-v2. The table has a dark background with light-colored text and icons. At the bottom, there is a "Citation" button.

T	Model	Average	ARC	HellaSwag	MMLU	TruthfulQA
◆	ICBU-NPU/FashionGPT-70B-V1.2	74.11	73.94	88.15	79.11	65.15
◆	sequelbox/StellarBright	74.1	72.95	87.82	71.17	64.46
◆	Riiid/sheep-duck-llama-2-70b-v1.1	74.07	73.04	87.81	70.84	64.58
◆	AIDC-ai-business/Marcoroni-70B-v1	74.06	73.55	87.62	70.67	64.41
◆	ICBU-NPU/FashionGPT-70B-V1.1	74.05	71.76	88.2	70.99	65.26
◆	adonlee/LLaMA_2_70B_LoRA	73.9	72.7	87.55	70.84	64.52
◆	uni-tianyan/Uni-TianYan	73.81	72.1	87.4	69.91	65.81
◆	Riiid/sheep-duck-llama-2	73.69	72.35	87.78	70.82	63.8
◆	Riiid/sheep-duck-llama-2	73.67	72.27	87.78	70.81	63.8
◆	fangloveskari/ORCA_LLaMA_70B_QLoRA	73.4	72.27	87.74	70.23	63.37
◆	ICBU-NPU/FashionGPT-70B-V3	73.26	71.08	87.32	70.7	63.92
◆	oh-yeontaek/llama-2-70B-LoRA-assemble-v2	73.22	71.84	86.89	69.37	64.79

HuggingFace OpenLLM Leaderboard

Key Metrics for LLM Evaluation

Evaluating large language models (LLMs) requires more than a simple pass/fail grade. Here are key metrics that show different parts of an LLM's abilities:

1. Accuracy and Facts

- Question Answering Accuracy: How well does the LLM answer questions based on facts? (e.g. SQuAD)

- Fact-Checking: Can the LLM identify and confirm factual claims in the text?

2. Fluency and Coherence

- BLEU/ROUGE Scores: Do the LLM's texts match up to human references in grammar and readability?
- Human Readability Score: How natural and well-organized are the LLM's texts, as judged by people?

3. Diversity and Creativity

- Unique Responses Generated: How many different responses can the LLM produce for a single prompt?
- Human Originality Score: How unique and unexpected are the LLM's texts, as judged by people?

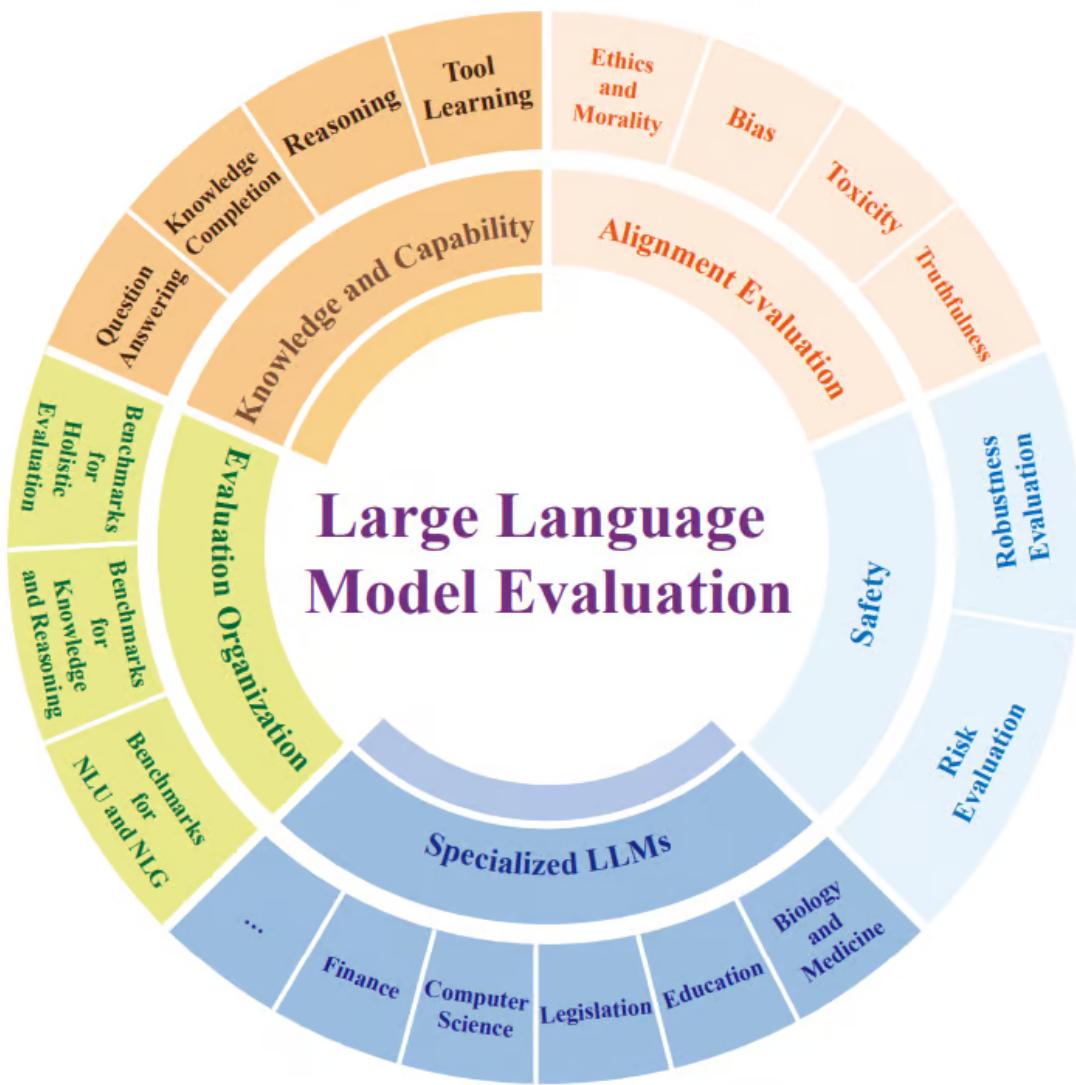
4. Reasoning and Understanding

- Natural Language Inference: How well can the LLM understand relationships between sentences? (e.g. MNLI)
- Causal Reasoning: Can the LLM make logical inferences and see cause-and-effect connections?

5. Safety and Robustness

- Resistance to Attack: How easily can cleverly crafted inputs mislead the LLM?
- Toxicity Detection: Can the LLM avoid generating harmful or offensive language?

No single metric gives the full picture. Use a balanced mix of metrics and human judgment to truly understand an LLM's strengths and weaknesses. This allows us to unlock their potential while ensuring responsible development.



Taxonomy of major categories and sub-categories of LLM evaluation (from this paper).

Future Directions in LLM Evaluation

While existing techniques like benchmark datasets, automatic metrics, and human evaluation offer valuable insights, the future of LLM evaluation needs a more comprehensive compass. We must chart a course towards:

- 1. Value Alignment and Dynamic Adaptation:** Evaluation should move beyond technical prowess and prioritize alignment with human values like fairness, explainability, and responsible text generation. Dynamic benchmarks that adapt to the ever-evolving nature of LLMs and real-world scenarios are crucial.
- 2. Agent-Centric and Enhancement-Oriented Measures:** Instead of isolated tasks, we need to evaluate LLMs as complete agents, assessing their ability to learn, adapt, and interact meaningfully. Moreover, evaluation should not solely identify flaws but offer metrics that guide improvement and suggest pathways for enhancement.

This future demands collaborative efforts from researchers, developers, and ethicists to

create evaluation methodologies that are not just rigorous, but also comprehensive and aligned with societal values. Imagine continuous testing in dynamic environments, evaluation metrics that prioritize fairness and responsibility, and collaborative efforts to ensure LLMs serve humanity in ethical and transformative ways. The journey toward truly meaningful LLM evaluation has just begun, and the future holds exciting possibilities for shaping the potential of these powerful language models.

Enjoy the weekend folks,

Armand

Whenever you're ready, there are FREE 2 ways to learn more about AI with me:

1. [**The 15-day Generative AI course**](#): Join my 15-day Generative AI email course, and learn with **just 5 minutes a day**. You'll receive concise daily lessons focused on practical business applications. It is perfect for quickly learning and applying core AI concepts. 10,000+ Business Professionals are already learning with it.
2. [**The AI Bootcamp**](#): For those looking to go deeper, join a full bootcamp with 50+ videos, 15 practice exercises, and a community to all learn together.



Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beehiiv

Subject: A Guide to Retrieval Augmented Generation
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Saturday, January 27, 2024 5:31:09 PM GMT+05:30
Date Received: Saturday, January 27, 2024 5:31:11 PM GMT+05:30

January 27, 2024 | [Read Online](#)

A Guide to Retrieval Augmented Generation

Learn about most popular Generative AI use case for Business



Welcome to the 1,154 new members this week! **nocode.ai** now has 35,549 subscribers

Retrieval-Augmented Generation (RAG) represents a significant evolution in the field of AI, particularly in enhancing the capabilities of large language models (LLMs). In today's guide, I'll provide a deep dive into the RAG essentials, exploring its functionality, importance, operational mechanisms, and best practices for implementation.

Today I'll cover:

- What is exactly RAG
- How it works
- Business use cases for RAG
- The advantages of utilizing RAG
- Recommended architectural framework
- RAG vs Fine-Tune
- A step-by-step implementation tutorial

Let's Dive In!

What is Retrieval Augmented Generation (RAG)?

Retrieval-augmented generation is an AI framework that integrates LLMs with external knowledge sources. This synergy allows LLMs to supplement their internal information with facts retrieved from an external knowledge base, such as customer records, dialogue paragraphs, product specifications, and even audio content. By utilizing this external context, LLMs can generate more informed and accurate responses.

Human work will shift from performing numerous manual look-ups and information gathering to directing teams of Large Language Models (LLMs) and synthesizing the results.

Enterprises will employ 100-1000s of these AI assistants across every job function.

Why is RAG Important?

RAG addresses several critical challenges faced by LLMs:

- **Knowledge Cutoff:** LLMs are limited by their training data. RAG provides access to external, up-to-date knowledge, enabling more accurate responses.
- **Hallucination Risks:** LLMs may produce factually inaccurate responses. RAG supplements LLMs with external information, reducing such risks.
- **Contextual Limitations:** LLMs often lack context from private data, leading to inaccuracies. RAG provides relevant, domain-specific data, ensuring more informed responses.
- **Auditability:** By citing sources, RAG improves the auditability of responses generated by GenAI applications.

How RAG Works

RAG operates in two distinct phases:

1. **Retrieval Phase:** Algorithms search and retrieve snippets of information relevant to the user's prompt from various data sources like document repositories, databases, or APIs.
2. **Content Generation Phase:** The retrieved context is provided as input to a generator model, typically an LLM. This model uses the context to generate text output grounded in relevant facts and knowledge.



[nocode.ai](#)

Diagram of How RAG Works

To ensure compatibility, document collections, and user queries are converted to numerical representations using embedding language models. The RAG architecture then compares these embeddings, augmenting the original user prompt with relevant context before sending it to the foundation model.

Business use cases for RAG

RAG technology has found many practical applications across various business sectors:

- **Enhanced Search Outcomes:** In healthcare, RAG improves the examination of Electronic Medical Records (EMRs) and the findings of clinical trials.
- **Interactive Data Conversations:** RAG simplifies complex data interactions, allowing non-technical stakeholders to query databases directly using natural language.
- **Advanced Customer Support Chatbots:** In industries like IT and manufacturing, RAG-equipped chatbots provide precise responses to customer inquiries.
- **Summarization for Efficiency:** In education, RAG streamlines the process of grading essays or creating condensed study materials.
- **Data-Driven Decision Making:** In finance and legal sectors, RAG assists in drafting contracts and condensing regulatory documents.

The advantages of utilizing RAG

There are many benefits of employing RAG in business contexts. First and foremost, it significantly enhances the accuracy and relevance of AI-generated content. This leads to improved customer satisfaction, as responses are more tailored and informative. Furthermore, RAG can streamline operations, reducing the time and resources spent on content generation and data retrieval tasks.

Advantages

- **Cost-effective Training:** RAG requires less computational power and data compared to extensive fine-tuning or training LLM processes.
- **Access to Various Knowledge Sources:** RAG combines internal knowledge with that from external databases, resulting in more accurate answers.
- **Enhanced Scalability:** RAG can handle large datasets and intricate queries, surpassing conventional LLMs limited by their context window size.

Challenges

- **Risk of Hallucinations:** RAG can still make errors if the database lacks certain information.
- **Managing Scalability:** Increasing the database size can complicate quick and efficient data retrieval.
- **Potential Biases:** Biases in the retrieval database can influence the responses, necessitating tools to identify and mitigate these biases.

The Generative AI Maturity curve

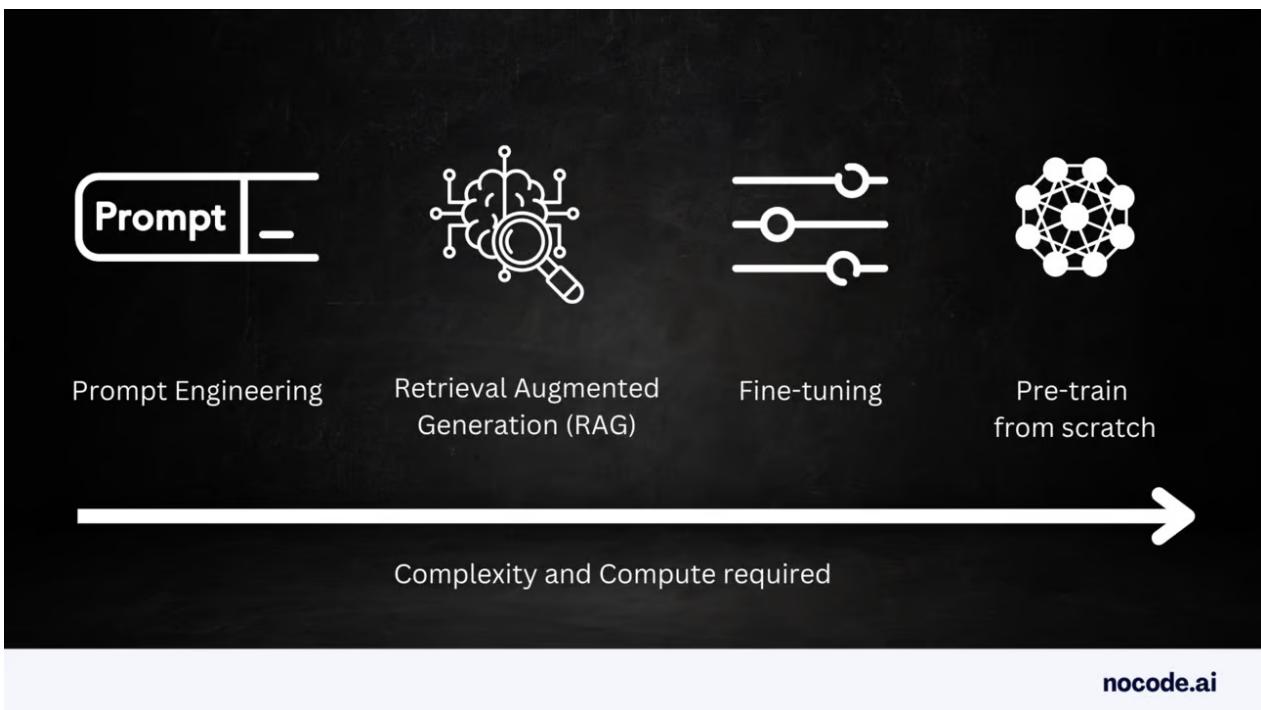
Before LLMs, ML development was linear. Teams needed months of tedious data collection, feature engineering, and multiple training runs, as well as a team of PhDs before the system could be produced as a customer-facing end product.

The table below summarizes key methods in LLM development, outlining their definitions, use cases, data requirements, advantages, and considerations for your business applications. It covers Prompt Engineering, Retrieval Augmented Generation (RAG), Fine-tuning, and Pretraining.

Method	Definition	Primary Use Case	Data Requirements	Advantages	Considerations
Prompt Engineering	Crafting specialized prompts to guide LLM behavior	Quick, on-the-fly model guidance	None	Fast, cost-effective, no training required	Less control than fine-tuning
Retrieval Augmented Generation (RAG)	Combining an LLM with external knowledge retrieval	Dynamic datasets and external knowledge	External knowledge base or database (e.g., vector database)	Dynamically updated context, enhanced accuracy	Increases prompt length and inference computation
Fine-tuning	Adapting a pretrained LLM to specific datasets or domains	Domain or task specialization	Thousands of domain-specific or instruction examples	Granular control, high specialization	Requires labeled data, computational cost
Pretraining	Training an LLM from scratch	Unique tasks or domain-specific corpora	Large datasets (billions to trillions of tokens)	Maximum control, tailored for specific needs	Extremely resource-intensive

Methods to implement LLMs in Business

Each of these techniques shows different stages of complexity and compute-intensiveness. What's cool about RAG is that it is very cost-efficient, and when implemented properly, the results can be excellent.



Generative AI Maturity curve and methods

Recommended Architectural Framework

Adopting RAG requires a thoughtful architectural approach. The blueprint suggests a framework that seamlessly integrates the retrieval and generative components. This includes robust databases, efficient indexing mechanisms for quick data retrieval, and a generative model that can effectively utilize the retrieved data. Ensuring smooth interoperability between these components is key to harnessing the full potential of RAG.

Components required in a RAG architecture:

- 1. Knowledge Base:** Think of this as RAG's library, filled with all sorts of information from documents, databases, or even APIs. It's like a treasure trove of knowledge for RAG to use.
- 2. User Query:** This is where you come in. You ask a question or make a request, and RAG starts its magic.

3. Retrieval Model:

- 1. Embedding Model:** This part turns the text from your question and the information in the knowledge base into numbers. It's like translating languages, but instead, it translates words into a form that the system can understand and compare.
- 2. Search Engine:** Armed with these numerical translations, RAG then searches through its library to find the most relevant information. It's like having a super-efficient librarian at your service.

4. Generation Model:

- 1. Large Language Model (LLM):** This is where RAG gets creative. It uses advanced text generation models (think of them as super-smart writing tools) like GPT-3 to craft a response that's both informative and easy to understand.

5. Integration and Orchestration:

- 1. Prompt Engineering:** This is a bit like scriptwriting for RAG. It takes the information found and mixes it with your original question to set the scene for the LLM.
- 2. Model Serving:** This is the backstage crew, making sure RAG gets your question and sends back the right answer.

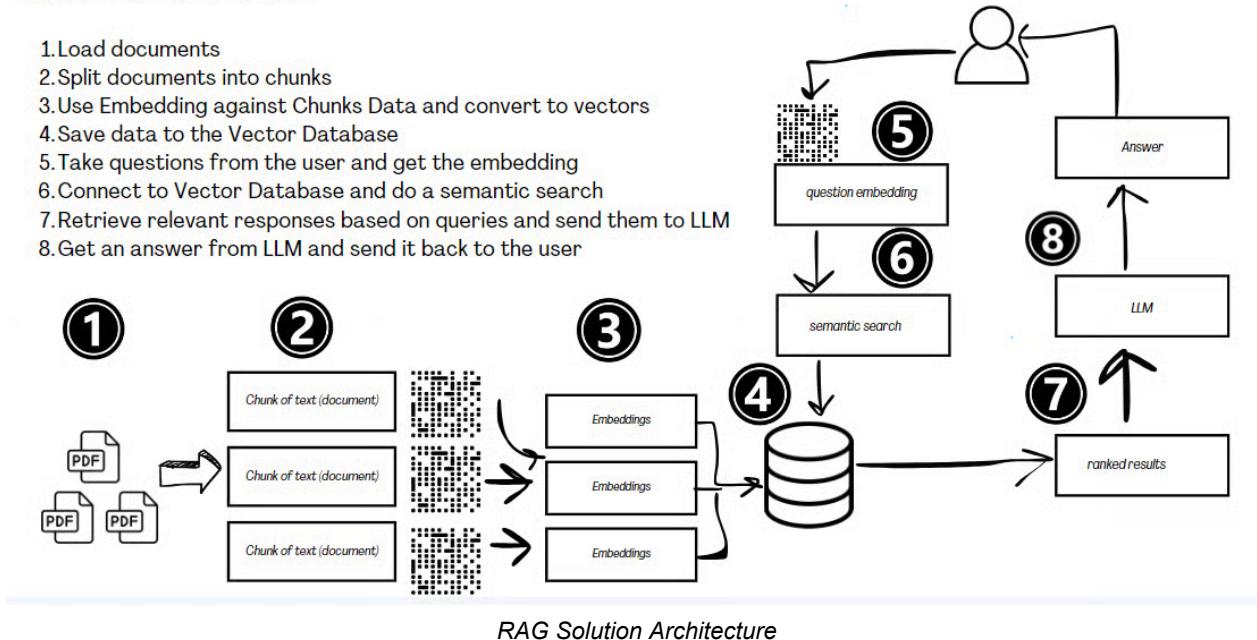
6. Extra Bits:

- 1. Monitoring and Logging:** Keeps an eye on RAG to make sure it's doing its job right.
- 2. User Interface:** This is where you interact with RAG, like in a chatbot or search engine.

A key component of the architecture is the Vector Database. It is used to store high-dimensional embeddings of text documents. Its primary role is to enable fast and efficient retrieval of information that is semantically similar to a given query. This retrieval is crucial for the RAG model to generate accurate and contextually relevant responses. The vector database ensures scalability, speed, and continuous updating of information, enhancing the overall performance of the RAG system.

Solution Architecture

1. Load documents
2. Split documents into chunks
3. Use Embedding against Chunks Data and convert to vectors
4. Save data to the Vector Database
5. Take questions from the user and get the embedding
6. Connect to Vector Database and do a semantic search
7. Retrieve relevant responses based on queries and send them to LLM
8. Get an answer from LLM and send it back to the user



The Grand Dilemma: RAG vs Fine-Tune

The debate in Generative AI often revolves around choosing between RAG and fine-tuning LLMs. This choice is influenced by the need for domain-specificity and the rate of data change.

I have put together a table to guide you through the decision-making process.

Aspect	Fine-Tuning	Retrieval-Augmented Generation (RAG)
Advantages	<ul style="list-style-type: none"> - Mitigates knowledge gaps with updated, specific data. 	<ul style="list-style-type: none"> - Provides near-real-time data updates.
	<ul style="list-style-type: none"> - Cost-effective compared to full model retraining. 	<ul style="list-style-type: none"> - Enhances transparency with source citations.
	<ul style="list-style-type: none"> - Suitable for training on private or specialized data. 	<ul style="list-style-type: none"> - Offers better data access control and personalization.
Challenges	<ul style="list-style-type: none"> - Struggles with frequent data updates. 	<ul style="list-style-type: none"> - Relies heavily on the efficiency of the search system.
	<ul style="list-style-type: none"> - Lacks clear traceability to original data sources. 	<ul style="list-style-type: none"> - Limited in context size provided to LLMs.
	<ul style="list-style-type: none"> - Potential for inaccurate information. 	<ul style="list-style-type: none"> - Possible over-reliance may curb model creativity.
Application	<ul style="list-style-type: none"> - Best for stable, less sensitive data. 	<ul style="list-style-type: none"> - Ideal for scenarios requiring real-time data relevance and flexibility.

Blending fine-tuning and RAG could leverage their respective strengths, using fine-tuning for stable, less sensitive data, and RAG for real-time data relevance and flexibility. This combination could offer a comprehensive solution in advanced Generative AI applications.

From my hands-on experience, it's clear: around 60% of current use cases are swiftly embracing the RAG approach, marking a transformative shift in practical AI application.

A step-by-step RAG tutorial

Let's get hands-on with the keyboard and try a simple RAG implementation ourselves. So much valuable information is trapped in PDFs and documents. In today's tutorial, we will build a chat interface that allows direct interaction with information in PDFs. We will complete the following steps:

1. Load documents
2. Split documents into chunks
3. Use Embedding against Chunks Data and convert to vectors
4. Save data to the Vector Database
5. Take data (question) from the user and get the embedding
6. Connect to Vector Database and do a semantic search
7. Retrieve relevant responses based on user queries and send them to LLM
8. Get an answer from LLM and send it back to the user

Watch the video tutorial: [link](#)

chat with your documents



nocode.ai

Play the tutorial

Do you prefer step-by-step written instructions? You can access them here: [link](#)

Conclusion

RAG is transforming the landscape of AI by enabling LLMs to generate responses that are not only accurate but also rich in context and relevancy.

Its applications span various domains, from chatbots to customer service, providing solutions that are both reliable and informed.

By adopting RAG and adhering to best practices in its implementation, businesses and developers can significantly enhance the capabilities of their AI-driven applications, ensuring they remain at the forefront of technological innovation.

Enjoy the weekend folks,

Armand

Whenever you're ready, there are FREE 2 ways to learn more about AI with me:

1. [**The 15-day Generative AI course**](#): Join my 15-day Generative AI email course, and learn with **just 5 minutes a day**. You'll receive concise daily lessons focused on practical business applications. It is perfect for quickly learning and applying core AI concepts. 10,000+ Business Professionals are already learning with it.
2. [**The AI Bootcamp**](#): For those looking to go deeper, join a full bootcamp with 50+ videos, 15 practice exercises, and a community to all learn together.



Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beejiiv

Subject: Learning AI Backwards
From: Tobias Zwingmann <ai4bi@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Friday, January 26, 2024 7:35:39 PM GMT+05:30
Date Received: Friday, January 26, 2024 7:35:39 PM GMT+05:30

January 26, 2024 | [Read Online](#)



Learning AI vs. Doing AI

How getting hands-on quickly helps you succeed with AI



[Tobias Zwingmann](#)



Hi there,

This week, I spoke at two events. One was about AI in Industry and the other about AI in Social Media. Two very different topics, but they shared one big idea: AI isn't here just for experts, but for everyone.

The key question is: How do you prepare or "enable" people for this? What's not working is to send everyone back to classroom and make them study hard.

I think it's time to flip the script and spoil the end. Let people access and experience modern AI right away. Introduce more technical knowledge as needed, depending on the user and their intentions.

Let's jump in!

Talking about hands-on AI, here's a great opportunity coming up:



In this [upcoming hands-on workshop](#), I'll walk you through 20+ use cases for ChatGPT in data analytics. LIVE on February 7 & 8, 2024. Recording available when you register.
(Subscription required, but you can join with a free trial)

In Applied AI, Practice Beats Theory

Fun fact: Even the world's leading AI experts can't fully explain how AI brains like ChatGPT really work in detail. It's still an open research question, as OpenAI's Chief Scientist Ilya Sutskever explains in this video:



Chat with OpenAI CEO and Co-founder Sam Altman, and Chief Scientist Ilya Sutskever

So as a user, don't worry too much about the intricacies. You can leave that to an armada of AI researchers to figure out. But the good thing is, in the meantime, you can focus on what you can control best:

Find out if and how AI can help you get better at what you already do.

Honestly, I never met a single person who got successful with AI by just reading about it. Yet, many who started with no clue built amazing use cases by jumping into ChatGPT and figuring things out.

And you don't need an AI degree for that.

There's a great German term for this - "Handlungswissen", best translated as "know-how". It describes practical skills, the so-called "ability" that has been acquired through physical experience and practice.

The enabler: AI accessibility

If you want to learn by doing, then you actually need to be able to do the thing and fail safely.

For example, if I wanted to become a pilot, it's probably a bad idea to board a plane and "figure it out". Both the cost of failure and the cost of access to the technology are just too high.

Five years or so ago, the only feasible way to get "hands-on" with modern AI was to:

- Get data
- Label data
- Train a model
- Deploy the model
- Find out that your "AI" wasn't even all that good.

Especially in businesses, that means a high cost of failure because you kept a group of people busy for probably some months.

Today, this has changed fundamentally with multi-purpose AI systems available, as-a-service.

For free.

ChatGPT, Bard, or Bing are just the most popular examples. More niche applications like [Magnific AI](#), [Codeium](#), [DeepL Write](#), or [Tome](#) appear by the day. There's literally NO

reason why you shouldn't try these out (responsibly, of course - without a critical use case or sharing private data).

And with that, we have the chance to totally change the way we gain knowledge about AI.



Greg Brockman
@gdb

X

It is hard to describe how much you learn by actually doing — by carefully considering all factors, making a decision, and then taking responsibility for the outcome. Unlocks wisdom that cannot be arrived at any other way.

Jan 23, 2024

5.41K Likes 785 Retweets 133 Replies

"Learning AI Backwards"

Typically, the first rides in an AI services like ChatGPT will be quite exciting. However, this effect will wear off quickly if you're not quite sure how to make the most of it. Chances are, you'll hit a wall pretty soon where AI answers are kind of "meh" or totally off track.

A lot of people's AI journey ends right here.

So, you might wonder: Is diving into AI hands-first really the right way to learn it? I would put it like this: Jumping straight into using ChatGPT is **the best first step**.

Let's flip the script and approach AI from the end, not the start.

Here's a 3 day action plan:

Day 1: Do something simple

Get ChatGPT Plus and fire up ChatGPT4. But don't sweat the complex stuff yet. Start simple. For example, let it write a recipe or have a critical discussion with you. Go play and have fun - but keep expectations low.

Day 2: Do something more complex that's worked for others

Before you get too advanced, see what other people have done successfully with ChatGPT. For example, check out this [data analytics prompt flow](#), or [build a simple assistant](#). Try to get more comfortable with the tool and push your own limits a little more.

Day 3: Take a break and reflect

Now, before you jump to any conclusions or enroll in a course, take a moment to reflect. Answer the question: "Could this potentially make a difference in my daily life or work?"

If it's a "no", no stress. Check back soon. But if it's a "yes" or a "maybe", that's your cue to dig a little deeper.

From AI Experience to AI Understanding

If you're still here, it's time to peel back the layers to [understand a little more about the basics](#), the do's and don'ts, and what AI might bring to the table in the future.

So, where do you start? Here are a few breadcrumbs to follow:

- **Intro:** [Elements of AI](#) – One of the world's most popular AI courses. Non-technical, up-to-date and free by the University of Helsinki.
- **Bootcamp:** [Nocode AI Bootcamp](#) – Taught by IBM's Director of AI, Armand Ruiz, this free, hands-on bootcamp follows the "learning by doing" approach over a 10-week course.
- **Comprehensive:** The [Complete Artificial Intelligence and ChatGPT Course](#) is an in-depth walkthrough of AI applications in business, with a special focus on ChatGPT, taught by AI expert Luka Anicin and business expert Chris Haroun.

If you want to dive even deeper, or build professional AI applications yourself, consider looking at the excellent Machine Learning / AI specializations by [Google](#), [Microsoft](#), or [Deeplearning.AI](#) - most of them are free.

Conclusion

Like any skill, AI proficiency comes from a mix of learning and doing. Today's AI systems make it easier than ever to get your hands on it. Embrace that, but don't stop there. "Experience" is only one component of learning, the other is "understanding the why."

It's about understanding the reasons behind AI's actions, its limitations, and its potential

impact on our world. These deeper insights will transform you from a naive user to an informed participant in the AI age.

Think of it like driving a car. Knowing how to operate the vehicle is essential, but understanding the rules of the road and why they exist makes you a skilled and safe driver.

Similarly, in AI, the "how" allows you to use the technology, but the "why" enables you to apply it wisely and creatively in different scenarios.

Stay curious, keep exploring.

See you next Friday!

Tobias

PS: If you found this newsletter useful, [please leave a testimonial!](#) It means a lot to me!



Update your email preferences or unsubscribe [here](#)

© 2024 AI For Business Growth

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beejiiv

Subject: Generative AI Course - Your Voice Matters !
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Friday, January 26, 2024 1:36:58 AM GMT+05:30
Date Received: Friday, January 26, 2024 1:36:59 AM GMT+05:30

January 25, 2024 | [Read Online](#)

New Post



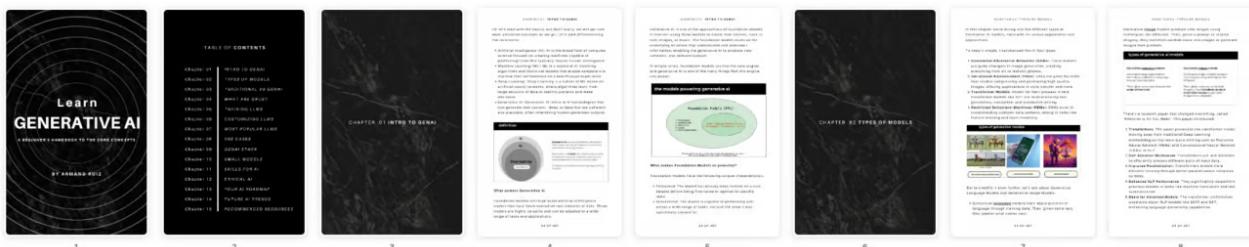
Hi!

What a journey it's been! As our 15-day Generative AI course concludes, I'd be thrilled if you could share [a quick testimonial about your experience](#). Your insights can inspire others and help me improve.

You can share your testimonial here:

[Share testimonial](#)

As a thank you, you'll receive [early access to my book](#) with all the lessons from this course. I will have it ready for you in February:



Learn Generative AI - Book coming soon

Thanks for being part of this adventure!

Warm regards,

Armand



Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beehiiv

Subject: The Generative AI Development Stack
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Saturday, January 20, 2024 5:32:30 PM GMT+05:30
Date Received: Saturday, January 20, 2024 5:32:31 PM GMT+05:30

January 20, 2024 | [Read Online](#)

The Generative AI Development Stack

Key elements to build tailored solutions



Welcome to the 2,194 new members this week! **nocode.ai** now has 34,483 subscribers

To build GenAI apps using your data, you need a modular, well-integrated stack that powers tailored solutions.

Today I'll cover:

- The Modern AI Stack
- Compute and Foundation Models
- Data: The Fuel for AI's Engine
- Deployment: Bringing AI to Life
- Observability: Watchtower of AI Systems
- RAG: The Architectural North Star for Today's AI
- The Future AI Revolution is still evolving

Let's Dive In!

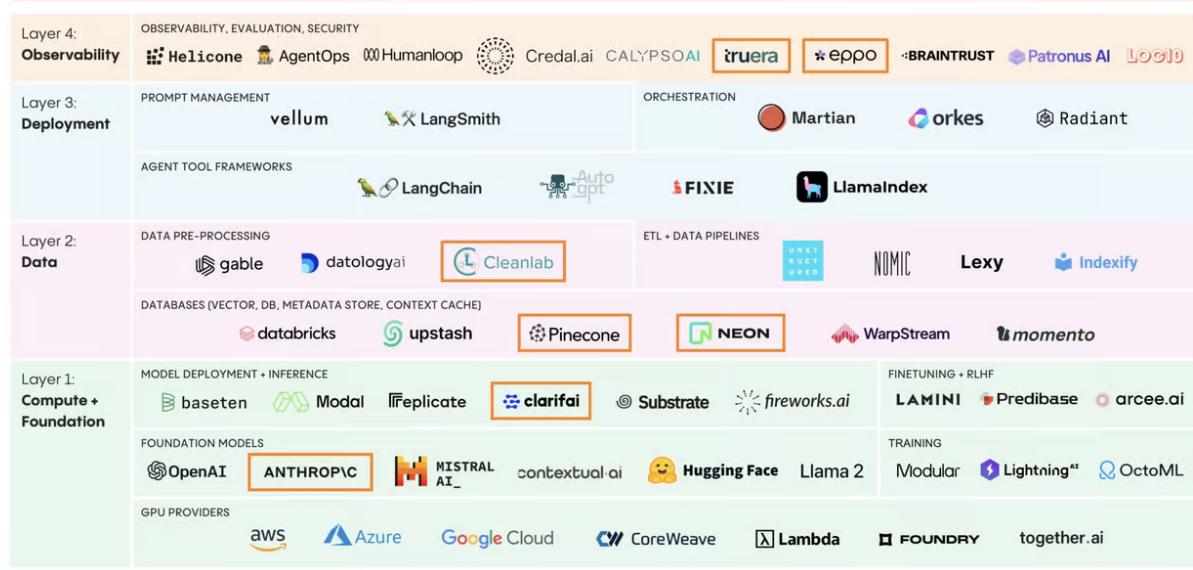
The Modern AI Stack

In 2023, the modern AI stack solidified its position as a generative AI powerhouse, with

record-high investments in the industry. [Menlo Ventures](#), identified four critical layers in this evolving AI landscape that I agree with:

1. **Compute and Foundation Models:** This foundational layer includes both the AI models and the infrastructure needed to train, optimize, and deploy them. It's the bedrock of the stack, enabling the core functionalities of AI applications.
2. **Data:** Serving as the lifeblood of AI, this layer focuses on the infrastructure that integrates [Large Language Models](#) (LLMs) with enterprise data systems. It encompasses data preprocessing, ETL (Extract, Transform, Load) processes, and various databases, ensuring that AI models have access to relevant and structured data.
3. **Deployment:** This layer is all about bringing AI models into practical use. It includes tools that help developers manage AI applications, ensuring smooth integration and consistent performance.
4. **Observability:** The final layer of the stack is dedicated to monitoring and security. It includes solutions that track the behavior of LLMs in real time and guard against potential threats, ensuring safe and reliable AI operations.

Modern AI Stack: The Emerging Building Blocks for GenAI



© 2024 Menlo Ventures

Backed by Menlo Ventures

Modern AI Stack from Menlo Venture

Compute and Foundation Models

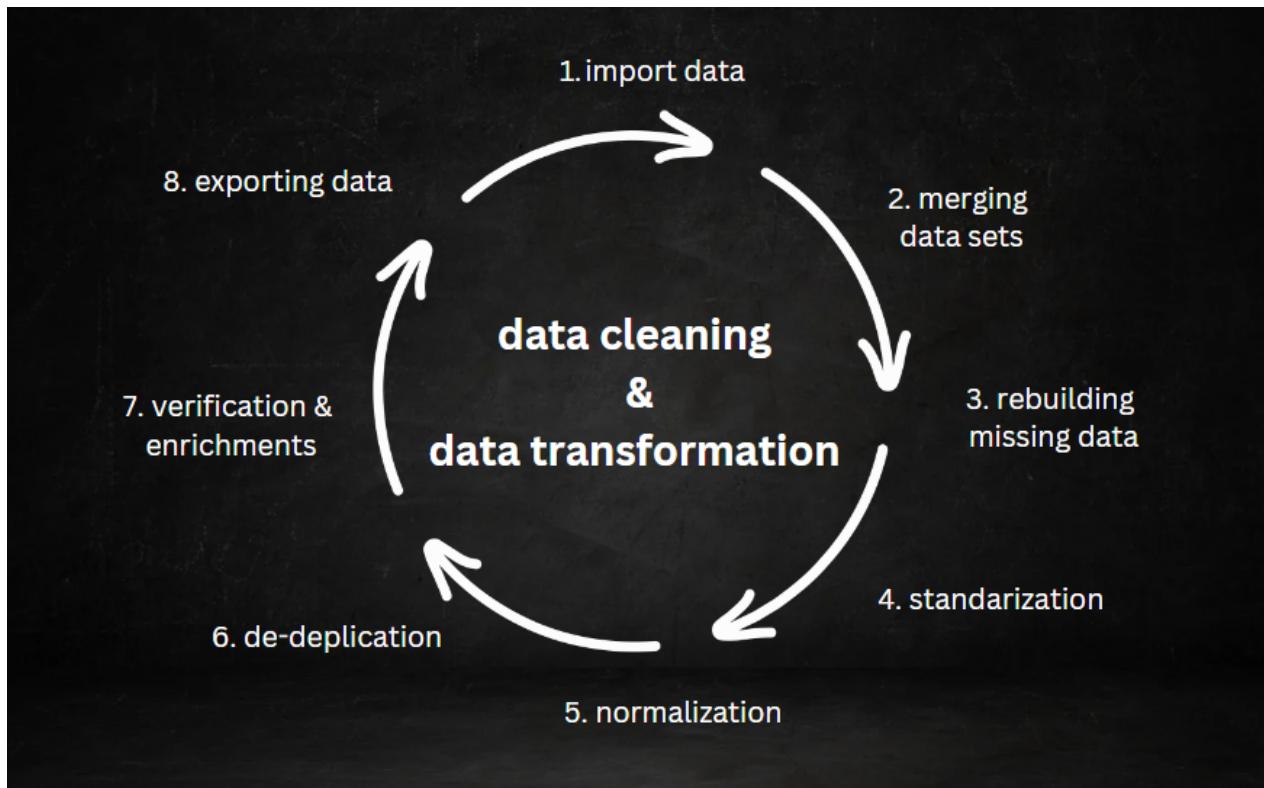
The compute and foundation model layer encompasses the core AI models and the infrastructure required for their training, fine-tuning, optimization, and deployment. Let's go through each of them:

1. **GPU Providers:** GPUs are crucial for processing AI algorithms. They handle complex computations required for model training and inference, with leading providers like NVIDIA and AMD at the forefront.
 2. **Foundation Models:** These are extensive, pre-trained models (like GPT, BERT) that form the basis for specialized AI applications, providing a broad understanding that can be tailored to specific tasks.
 3. **Training:** Training involves using large datasets to educate models, enabling them to make informed predictions and decisions. It's a complex process, balancing data diversity and computational demands.
 4. **Fine-tuning + RLHF:** Fine-tuning adjusts these models to specific needs, while Reinforcement Learning from Human Feedback (RLHF) further refines them based on human interactions, aligning AI outputs with real-world requirements.
 5. **Model Deployment:** The final step is deploying these models into applications, a process that involves strategic considerations like cloud platforms, edge computing, and scalability.
-

Data: The Fuel for AI's Engine

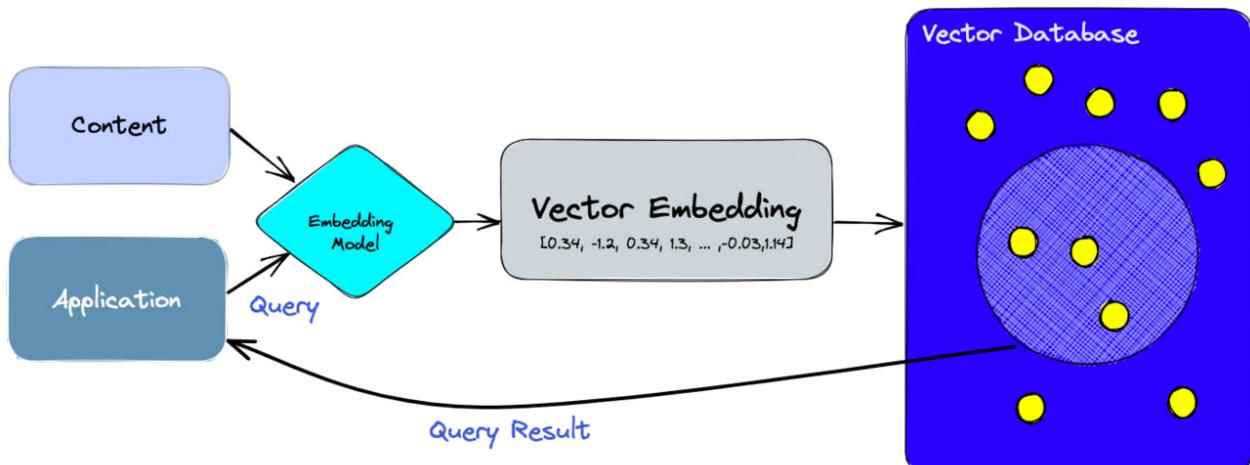
Data is the lifeblood of AI systems. Let's examine how the Data layer acts as the fundamental fuel driving AI systems, focusing on the integration, processing, and utilization of data for AI applications.

- **Integration and Processing:** Essential for customizing AI to specific business needs through the effective collection and transformation of diverse data sets.
- **ETL Processes:** Fundamental in preparing data for AI use, these processes standardize data into a usable format for AI analysis.



Steps Required in Data Integration and Processing

- **Vector Databases:** Vital for handling complex data, converting complex data (like text or images) into numerical vectors, enabling efficient processing for tasks like similarity searches and pattern recognition, essential for handling large datasets in AI applications.



Converting content into Vectors and storing it in a Vector Database

Deployment: Bringing AI to Life

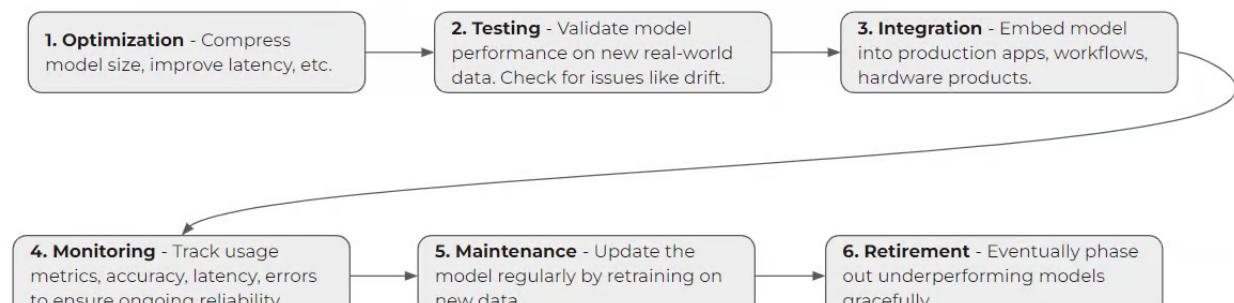
In the Deployment layer, we focus on crucial elements for launching AI models into real-

world applications:

- **Prompt Management:** Essential for fine-tuning AI responses to user inputs, ensuring relevant and accurate interactions.
- **Agent Tool Frameworks:** Provides tools for deploying AI agents, enabling them to perform tasks autonomously in varied environments.
- **AI Orchestration:** Involves coordinating multiple AI models and processes to ensure seamless, efficient operation, crucial for complex AI systems in enterprise settings.

This layer is key to transitioning AI from concept to practical application, aligning technical capabilities with user and business needs.

steps to deploy AI in production



Steps to deploy AI in production

Observability: Watchtower of AI Systems

The Observability layer in AI systems is crucial for ensuring transparency, security, and trustworthiness, while also facilitating performance evaluation.

- **Observability:** Essential for monitoring AI behavior and performance, ensuring transparency, and the ability to fine-tune systems for optimal operation.
- **Evaluation:** Involves assessing AI model performance and accuracy to guarantee they meet objectives and provide reliable results.

- **Security:** Focuses on protecting AI systems from threats like data breaches and model tampering, crucial for maintaining data integrity.
- **Trustworthiness:** Addresses the need for AI systems to be fair, ethical, and unbiased, building confidence in their use and deployment.

Overall, this layer is critical for maintaining the efficiency, security, and ethical standards of AI systems, ensuring their ethical application.

RAG: The Architectural North Star for Today's AI

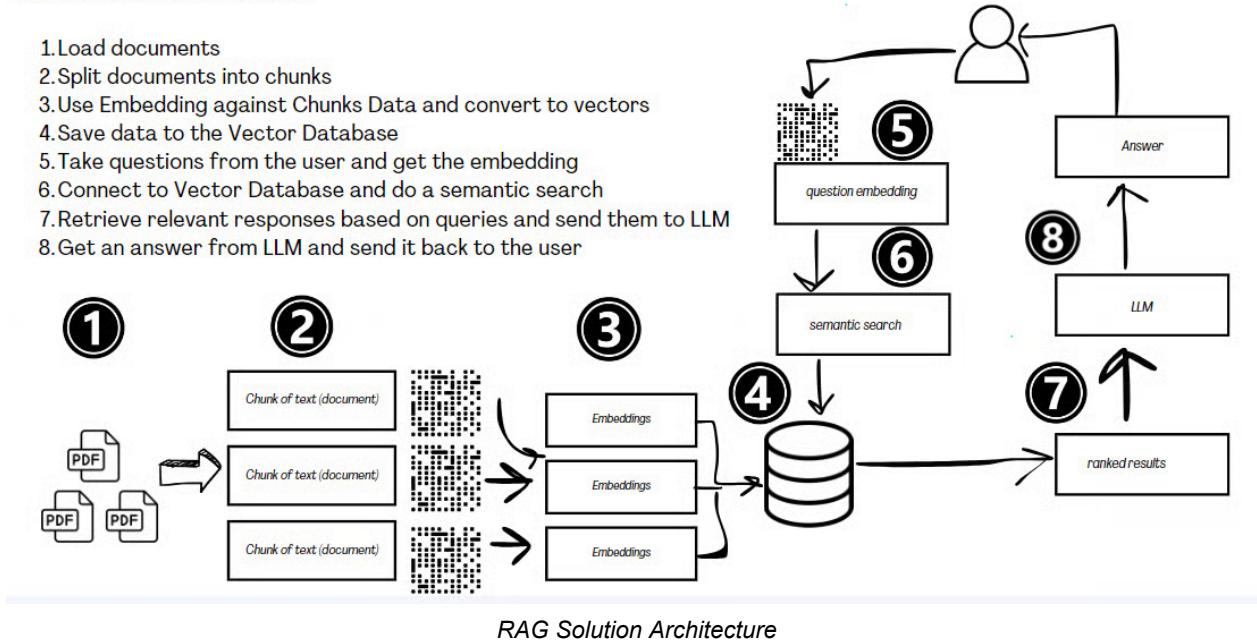
RAG (Retrieval Augmented Generation) enhances LLMs by integrating database access, improving accuracy and contextual relevance in AI outputs. More scalable and cost-effective than traditional methods, RAG reduces errors in LLMs, increasing trustworthiness and adaptability across various industries.

- **RAG in Action:** Enhances various AI applications, such as chatbots and content generation, by integrating external knowledge for depth and relevance.
- **Advantages:**
 - **Contextual Awareness:** Provides access to vast information, making outputs more relevant.
 - **Flexibility and Scalability:** Adapts to different domains with dynamic data integration.
 - **Enhanced Accuracy:** Combines generative capabilities with external data for improved reliability.
- **Future of RAG:** Points towards AI systems that intelligently integrate diverse data sources, promising more sophisticated, responsive, and autonomous AI solutions.

In essence, RAG is a transformative element in AI architecture, essential for creating more intelligent and adaptable AI systems. See the architecture below as a simple, yet very powerful example of a blueprint to develop RAG for most use cases.

Solution Architecture

1. Load documents
2. Split documents into chunks
3. Use Embedding against Chunks Data and convert to vectors
4. Save data to the Vector Database
5. Take questions from the user and get the embedding
6. Connect to Vector Database and do a semantic search
7. Retrieve relevant responses based on queries and send them to LLM
8. Get an answer from LLM and send it back to the user



The Future AI Revolution is still evolving

Emerging developments shaping the future of AI include:

- **Advanced RAG Applications:** Addressing limitations in current RAG systems with sophisticated techniques for improved accuracy and context understanding.
- **Rise of Small Models:** Shifting focus to smaller, task-specific models for efficiency, supported by advanced ML pipeline infrastructure and quantization techniques.
- **Innovations in Observability and Evaluation:** Developing automated tools for better AI performance review and reliability, moving beyond manual evaluations.
- **Serverless Architectures in AI:** Transitioning towards serverless solutions for AI infrastructure to simplify operations and enhance agility and scalability.

I hope you enjoyed today's newsletter.

See you again next week!

Armand

Whenever you're ready, there are FREE 2 ways to learn more about AI with me:

1. [The 15-day Generative AI course](#): Join my 15-day Generative AI email course, and learn with just 5 minutes a day. You'll receive concise daily lessons focused on practical business applications. It is perfect for quickly learning and applying core AI concepts. 10,000+ Business Professionals are already learning with it.
2. [The AI Bootcamp](#): For those looking to go deeper, join a full bootcamp with 50+ videos, 15 practice exercises, and a community to all learn together.



Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beejiiv

Subject: Day 15: Continuing Your AI Journey
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Wednesday, January 17, 2024 2:47:57 PM GMT+05:30
Date Received: Wednesday, January 17, 2024 2:48:29 PM GMT+05:30

Day 15 - Continuing Your AI Journey

We've reached the final day of our 15-day exploration into Generative AI. It's not an end, but a beginning to a lifelong journey in AI learning and exploration. Here's how you can continue growing in this exciting field:

Online Courses & Tutorials

I created the AI Bootcamp that can provide you next level of practice. It includes videos and exercises with multiple tools, no coding knowledge is required. Get your free access here:

<https://www.nocode.ai/the-ai-bootcamp/>

If you would like the next level of depth, I like the Short Courses from DeepLearning.ai, they include topics such as:

- [Functions, Tools and Agents with LangChain](#)
- [Vector Databases: from Embeddings to Applications](#)
- [Quality and Safety for LLM Applications](#)
- [Building and Evaluating Advanced RAG Applications](#)
- [Reinforcement Learning from Human Feedback](#)

These require basic Python knowledge.

Follow leaders in the AI space

These are some of the AI leaders I follow on LinkedIn:

- [Aishwarya Srinivasan](#): Data Scientist | LinkedIn Top Voice Data & AI | EB1A Recipient | 460k+ Followers | Ex- Google, Ex-IBM
- [Allie K. Miller](#): AI Entrepreneur, Advisor, and Investor | 1MM+ followers | Former Amazon, IBM | LinkedIn Top Voice for AI 2019-2023
- [Luis Serrano](#): AI scientist | YouTuber - 120K followers | Author of Grokking Machine Learning
- [Andriy Burkov](#): ML at TalentNeuron, author of The Hundred-Page Machine Learning Book and the Machine Learning Engineering book

I also publish content every day with tips and best practices to apply AI for Business. Feel free to follow me [here](#)

Join Newsletters

Don't fall behind on AI. These newsletters summarize the daily news so you get the latest AI trends and tools you need to know. I read them all every day.

- [The Rundown](#)
 - [Ben's Bites](#)
 - [The Neuron](#)
-

Remember, the field of AI is ever-evolving, and continuous learning is key. Embrace curiosity, keep experimenting, and stay connected to the AI community. You're now equipped to take the next steps in your AI journey.

I hope you enjoyed the course, feel free to write me with feedback to keep making it even better

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

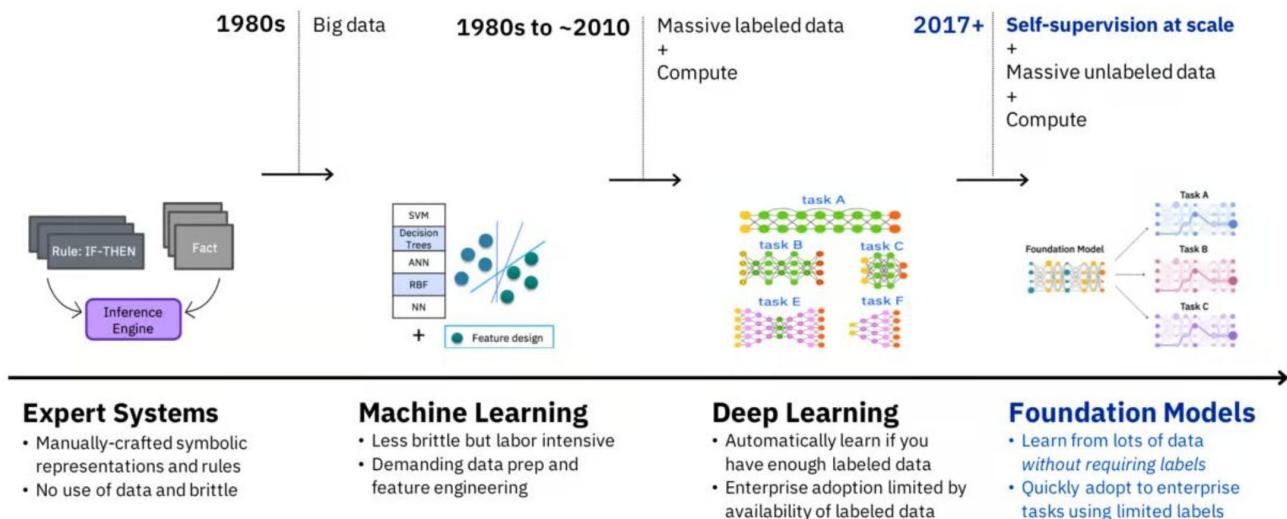
Subject: Day 14: Future Trends in AI
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Tuesday, January 16, 2024 2:47:20 PM GMT+05:30
Date Received: Tuesday, January 16, 2024 2:47:23 PM GMT+05:30

Day 14 - Future Trends in AI

Welcome to Day 14 of our 15-day journey through the world of Generative AI. Today, we're venturing into the future to explore the exciting innovations we can expect in the next decade. Let's dive in!

Where Do We Come From?

Understanding AI's journey is key to predicting its future. The past decade laid the foundation for advancements in machine learning and neural networks. Now, we're poised to build on this legacy, driving AI towards more sophisticated and nuanced applications.



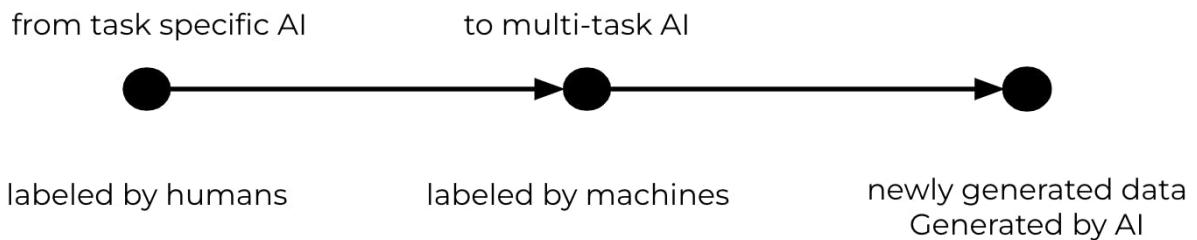
In **traditional machine learning**, individual siloed models require task-specific training and a significant amount of human-supervised learning. The limit on performance & capabilities for supervised learning are humans.

In contrast, **foundation models** are massive multi-tasking systems, adaptable with little or no additional training, utilizing pre-trained, self-supervised learning techniques. The limit on performance & capabilities is mostly on computing and data access (not labeling).

Synthetic Data

If the limit to a better model is more data, why don't create it artificially? The rise of synthetic data is a game-changer. It's about creating artificial datasets that can train AI without compromising privacy or relying on scarce real-world data. This innovation is set to revolutionize fields from healthcare to autonomous driving, making AI training more accessible, ethical, and comprehensive.

how do we get more data?



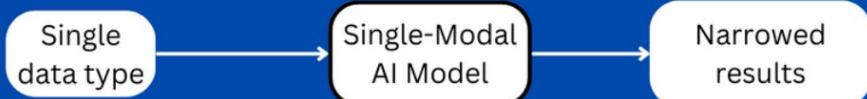
We will see increased usage of synthetic data because it overcomes data scarcity by allowing limitless generation, ensures balanced data distribution to avoid biases, and is cost-effective, bypassing the expensive and time-consuming process of real-world data collection and labeling.

Multimodality

Multimodality is the future of AI's interaction with the world. By integrating text, image, sound, and more, AI can understand and respond to complex queries with unprecedented accuracy. This holistic approach will deepen AI's integration into daily life, from smarter virtual assistants to more intuitive educational tools.

train ai with data generated by ai

Single-modal AI Model



Multimodal AI Model?



Reinforcement Learning

Want to take it to the next level? What if you could train the AI without data? Meet Reinforcement Learning, a technique poised to make significant strides. By learning through trial and error, AI systems will become more autonomous and capable of solving complex, real-world problems. This means smarter algorithms in everything from financial forecasting to climate change modeling.

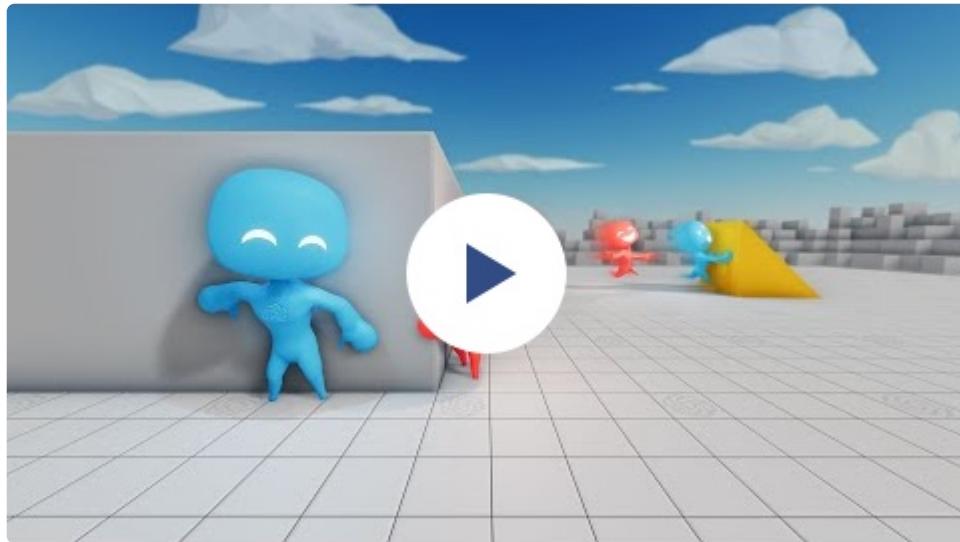
train AI without data

from task specific AI

to multi-task AI



Check out this video to see a demonstration of Reinforcement Learning in action:



Giving Tools to AI

In the next decade, AI will evolve beyond being just a tool to becoming a creator of tools. We will see AI designing software, crafting algorithms, and contributing to AI research. AI agents will autonomously manage projects in a continuous loop, executing tasks, enhancing results, and generating new tasks based on objectives and past outcomes. Their workflow will include task execution, result enrichment, task creation, and prioritization. Equipped with integration capabilities, these agents will be able to search for information in CRMs, access databases, send emails, and more. Frameworks like BabyAGI and Auto-GPT are already emerging to test these concepts.

Beyond Chats

While chatbots and conversational AI have made leaps, the future extends far beyond text. Expect AI that can seamlessly interact across various formats, offering richer, more immersive experiences. Whether it's in education, entertainment, or customer service, AI will engage us in more meaningful, dynamic ways.

As we approach the final day of our course, remember that the future of AI is not just about technology; it's about the creativity and ingenuity of those who wield it. Stay tuned for our concluding insights tomorrow!

See you tomorrow on our last day,

Armand

Subject: Day 13: Create Your AI for Business Roadmap
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Monday, January 15, 2024 2:47:47 PM GMT+05:30
Date Received: Monday, January 15, 2024 2:47:49 PM GMT+05:30

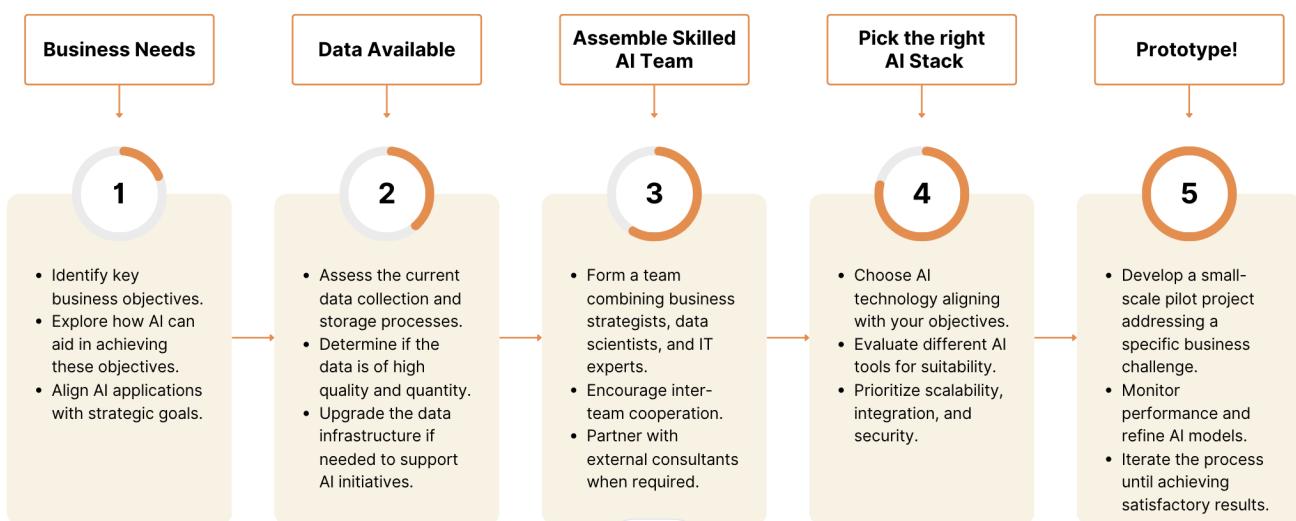
Day 13 - Create Your AI for Business Roadmap

Welcome to Day 13 of our AI in 15 Days email course. Today, we cover how to plan AI in business by creating an effective AI Business Roadmap. As leaders in your respective fields, the strategic integration of AI can be a transformative step for your organization, driving efficiency, innovation, and substantial growth.

AI Business Roadmap: A Strategic Necessity

An AI Business Roadmap isn't just a technical layout; it's a comprehensive plan that aligns AI technologies with your specific business goals. This blueprint encompasses timelines, resources, and technical requirements, and addresses potential risks and human factors ensuring a smooth AI integration.

AI Business Roadmap



Creating Your Roadmap: A Step-by-Step Guide

- 1. Identify Business Objectives:** Understand how AI can help achieve your goals, whether it's through automation, predictive analytics, AI chatbots, or innovative product development.

2. **Evaluate Data Infrastructure:** AI needs quality data. Assess your data collection, storage, and cleanliness to ensure your AI initiatives can thrive.
3. **Assemble a Skilled Team:** Combine business insight, technical skills, and data science. Include business strategists, AI specialists, and IT professionals, or seek external expertise as necessary.
4. **Choose Appropriate AI Technology:** Select AI tools like ML, NLP, RPA, or Computer Vision, aligned with your business needs.
5. **Prototype Development:** Start small with a pilot project to address specific challenges, refining AI models based on performance.
6. **Scale and Optimize:** Expand successful prototypes, integrating them into broader business operations and continuously optimizing.
7. **Implement Change Management:** Develop strategies to assist your workforce in adapting to AI, including training and understanding AI benefits.

High Impact, Low Effort AI Use Cases

Incorporate high-impact, low-effort AI applications such as AI chatbots for customer service, content generation, personalized marketing, predictive maintenance, routine task automation, fraud detection, demand forecasting, and streamlined recruitment. These use cases provide a strong foundation for your AI journey, ensuring meaningful returns with minimal initial effort.

Please share with me your roadmap plans for AI in 2024!

Best,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Day 12: Ethical Considerations in AI
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Sunday, January 14, 2024 2:47:18 PM GMT+05:30
Date Received: Sunday, January 14, 2024 2:47:20 PM GMT+05:30

Day 12 - Ethical Considerations in AI

As we approach the end of our AI series, Day 12 is dedicated to understanding the ethical considerations in AI, particularly the risks and responsibilities associated with deploying these powerful technologies.

AI is the most powerful technology ever created.

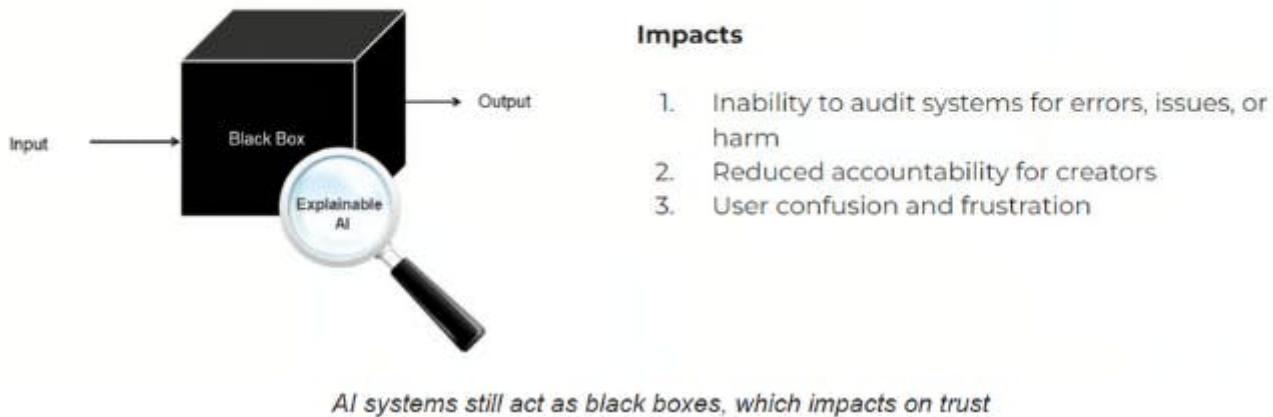
Understanding the Risks of AI

- **Bias, Fairness, and Accuracy:** AI systems can inadvertently replicate human biases present in training data, leading to unfair and inaccurate outcomes. Continuous improvement in data and training practices is crucial to enhance AI fairness.
- **Job Disruption:** AI's potential to automate tasks poses challenges for the job market, necessitating re-skilling initiatives and policy responses.
- **Weaponization:** The integration of AI into weapons systems raises ethical concerns about autonomy in warfare and the need for international governance.
- **Cybersecurity Vulnerabilities:** AI systems can be exploited for adversarial attacks, data poisoning, and model hacking, underscoring the importance of robust cybersecurity measures.
- **Misinformation Spread:** AI's ability to generate convincing fake news requires vigilance and tools to detect and mitigate the spread of misinformation.
- **Emergence of AGI:** The prospect of Artificial General Intelligence (AGI) adds another layer of ethical complexity, with implications for the future of humanity.

AI creating extinction risk for humanity is widely overhyped. AI develops gradually, and the “hard take off” scenario, where AI suddenly achieves superintelligence overnight is not realistic.

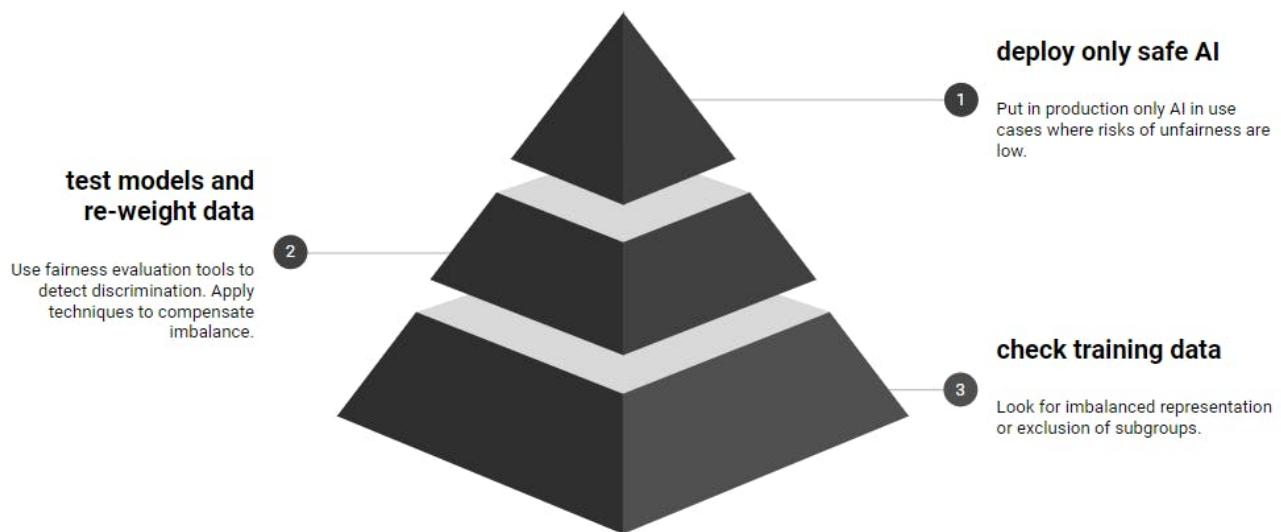
lack of transparency

Most AI systems act as "black boxes", with inner workings opaque to users.



The Role of AI Engineers in Mitigating Risks

- **Infrastructure Development:** AI Engineers are responsible for developing and managing AI infrastructure, ensuring systems are robust, scalable, and secure.
- **Ethical Implementation:** Part of their role involves applying ethical AI practices, such as prompt engineering and data management, to minimize risks like bias and inaccuracy.
- **Collaboration and Best Practices:** AI Engineers must collaborate across functions to promote AI best practices and ethical standards within their organizations.



Mitigating Harmful AI Outputs

- **Hallucinations and Fabrications:** AI systems can generate plausible but incorrect content. It's essential to design systems that minimize these risks.

- **Data Poisoning and Toxic Language:** AI must be safeguarded against harmful data inputs and language generation, requiring ongoing model training and content filtering.
- **Unstable Task Performance:** Addressing the inconsistent performance of AI models involves careful prompt engineering and verification processes.
- **Human Oversight:** Incorporating human review in high-stakes AI applications provides an additional layer of safety and accuracy.

The Path Forward: Ethical AI

- **Balancing Act:** The challenge lies in leveraging AI's potential while responsibly managing its risks and ethical implications.
- **Responsible Deployment:** Implementing safety checks, transparency measures, and human oversight is crucial for ethical AI deployment.
- **Ongoing Education and Policy Development:** Continuous learning and policy evolution are necessary to keep pace with AI advancements and ensure its beneficial use.

Foundation properties for AI Ethics

AI ethics provides guidelines to mitigate risks like amplifying human cognitive biases, enable responsible innovation, avoid harmful outcomes, and maintain public trust given AI's potential to rapidly scale both benefits and risks.

fairness	accountability	transparency	privacy & security

AI holds tremendous promise for transforming businesses and society. However, it's crucial to approach AI development and deployment with a keen awareness of its ethical implications, ensuring that these powerful tools are used responsibly and for the greater good.

Tomorrow more!

Best,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Why AI needs so much Energy
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Saturday, January 13, 2024 5:32:02 PM GMT+05:30
Date Received: Saturday, January 13, 2024 5:32:03 PM GMT+05:30

January 13, 2024 | [Read Online](#)

Why AI needs so much Energy

The Environmental Cost of AI Computational Power



Welcome to the 4,879 new members this week! **nocode.ai** now has 32,434 subscribers

AI has delivered amazing advances, but its hunger for data and computing power is also taking a toll on the environment. This raises important questions about AI's energy usage and carbon footprint that we need to address.

Today I'll cover:

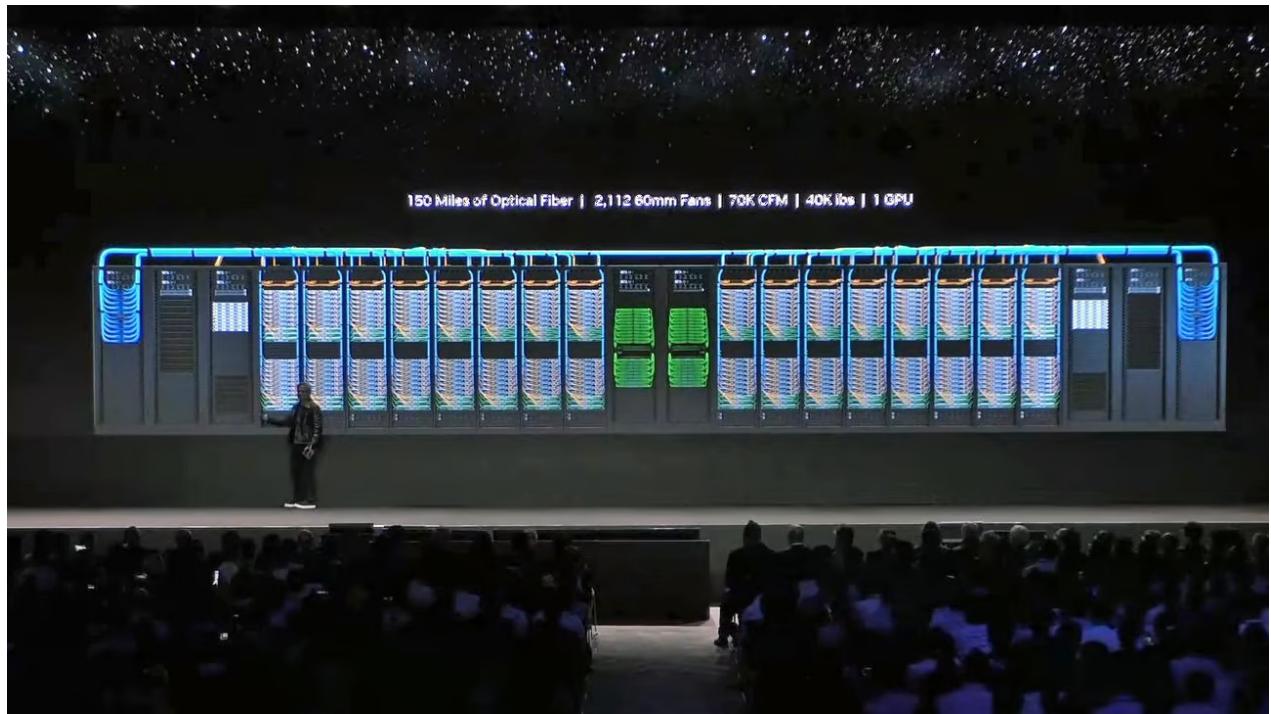
- Why AI consumes so much Energy
- Environmental Impact of High Energy Consumption
- Case Studies: AI's Energy Usage in the Real World
- Sustainable Solutions and Innovations
- Industry Initiatives and Corporate Responsibility
- Looking Ahead: The Future of AI and the Environment

Let's Dive In!

Why AI Needs So Much Energy

AI systems, especially machine learning models like deep neural networks, require vast

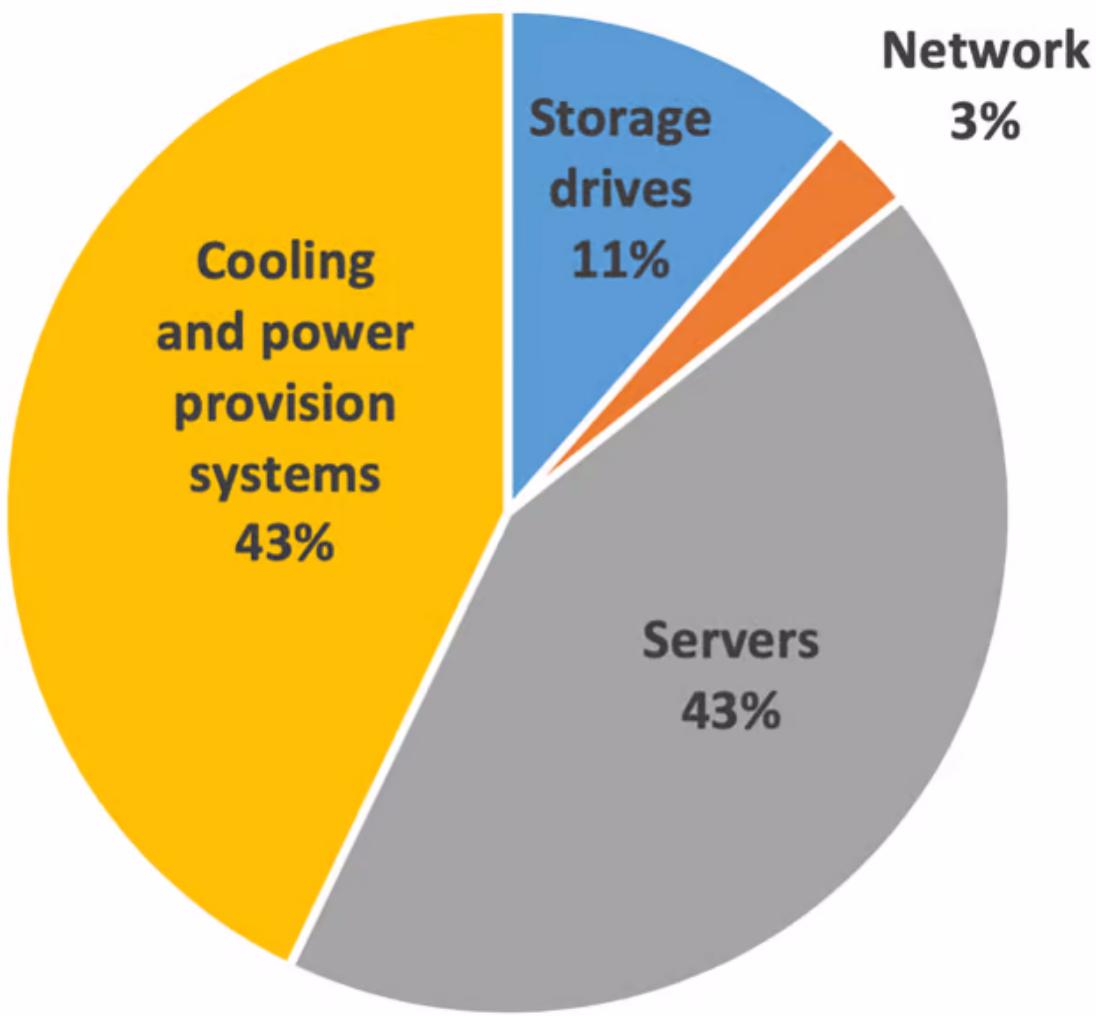
amounts of computing power during the training process. Training a single large AI model can emit as much carbon as five cars over their lifetimes. This intense computation relies on electricity from power grids powered by fossil fuels. As AI adoption grows, its carbon footprint expands.



Jensen Huang presenting AI Supercomputer DGX GH200

In an AI factory, more commonly known as a data center, energy consumption is spread across several key areas:

1. **Servers:** The core of a data center, they perform intense computing tasks for AI, significantly contributing to energy usage.
2. **Cooling Systems:** Essential for regulating the temperature of servers, these systems are major energy consumers.
3. **Storage Systems:** Used for housing vast amounts of data, storage systems require energy for operation and cooling.
4. **Networking Equipment:** Devices like routers and switches, necessary for connectivity, add to the overall energy consumption.
5. **Power Infrastructure:** Components like UPS and voltage regulators ensure steady power supply but also contribute to energy use.

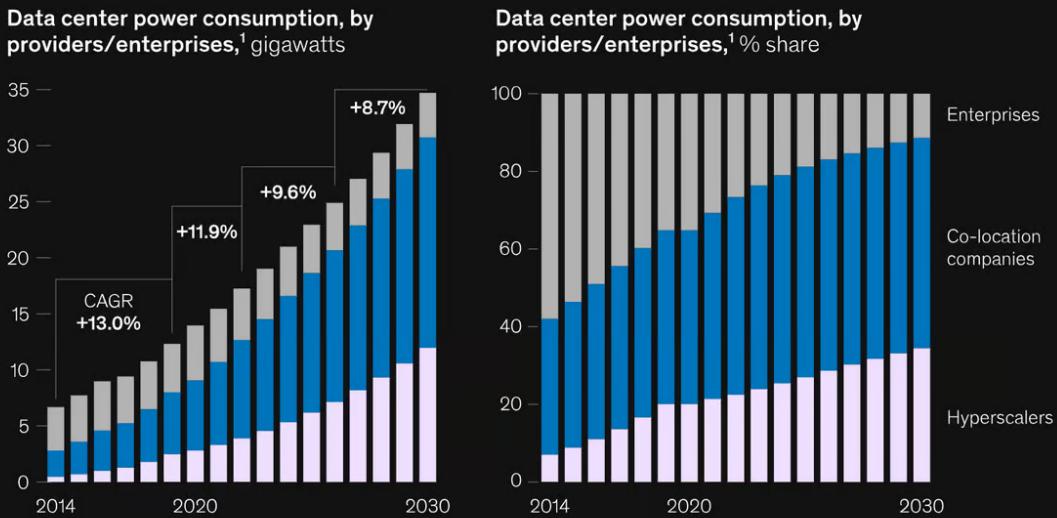


How much energy do data centers use? Image from Aspen Global Change Institute

Environmental Impact of High Energy Consumption

The environmental impact of AI's energy consumption is a growing concern. Data centers, the backbone of AI operations, are estimated to account for about 1% of global electricity use. This figure may seem small at first glance, but it's significant when considering the carbon footprint associated with electricity generation. Many data centers still rely on fossil fuels, contributing to greenhouse gas emissions. Also, the forecast shows 10% growth in data center demand every year.

US data center demand is forecast to grow by some 10 percent a year until 2030.



¹Demand is measured by power consumption to reflect the number of servers a data center can house. Demand includes megawatts for storage, servers, and networks.

McKinsey & Company

McKinsey & Company report of US data center demand forecast

The carbon footprint of AI is not limited to energy consumption during operations. The entire lifecycle of AI hardware, from manufacturing to disposal, has environmental implications. The production of specialized AI components, like GPUs, involves energy-intensive processes and the use of rare earth elements, which have their own environmental and ethical issues.

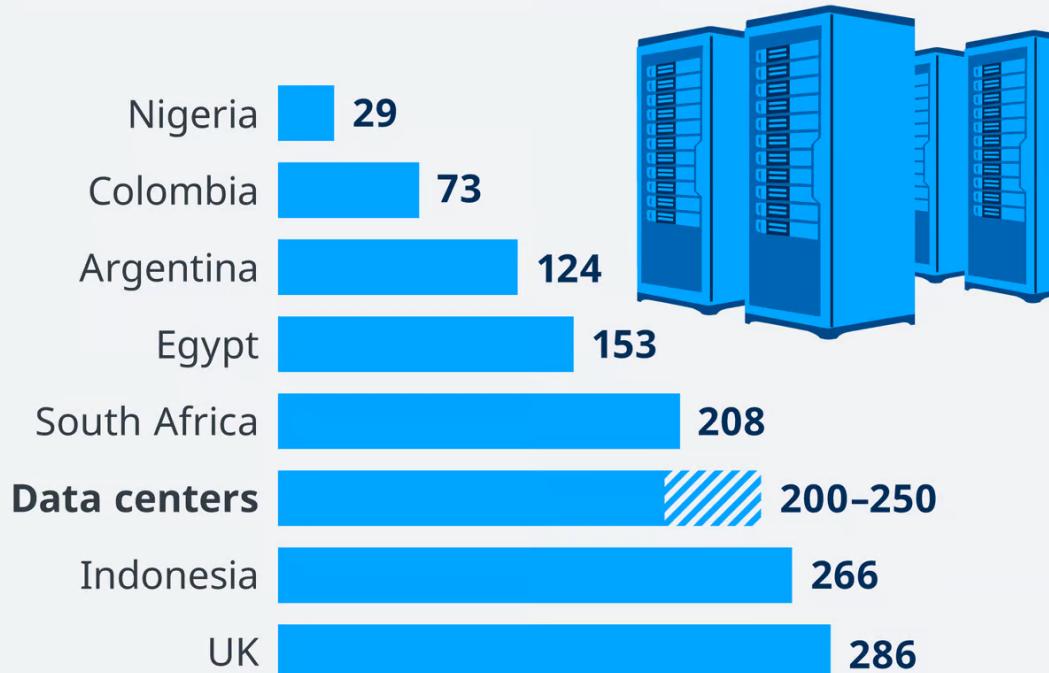
Case Studies: AI's Real-World Energy Usage

Here are some notable case studies

- 1. Google DeepMind's AI System (2021):** Consumed about 652,000 kWh yearly, nearly eight times more than the average UK household's usage, highlighting AI's high energy demands.
- 2. OpenAI's GPT-3 Model:** Required energy equivalent to burning 9 million pounds of coal for its pre-training, showcasing the significant carbon footprint of large AI models.
- 3. AI in Cloud Services (Google, Amazon):** Despite improving operational efficiency, AI optimizations in cloud services still entail substantial energy consumption, emphasizing the environmental cost of advanced AI.

Data centers use more electricity than entire countries

Domestic electricity consumption of selected countries vs. data centers in 2020 in TWh



Source: Enerdata, IEA



Sustainable Solutions and Innovations

How can the tech industry support AI's energy and computational demands sustainably?

Initiatives include improving computing efficiency, utilizing renewable energy sources like solar and wind, recycling waste heat into electricity, and developing carbon offset programs around AI system usage. Nuclear energy is also considered the best source given the following points:

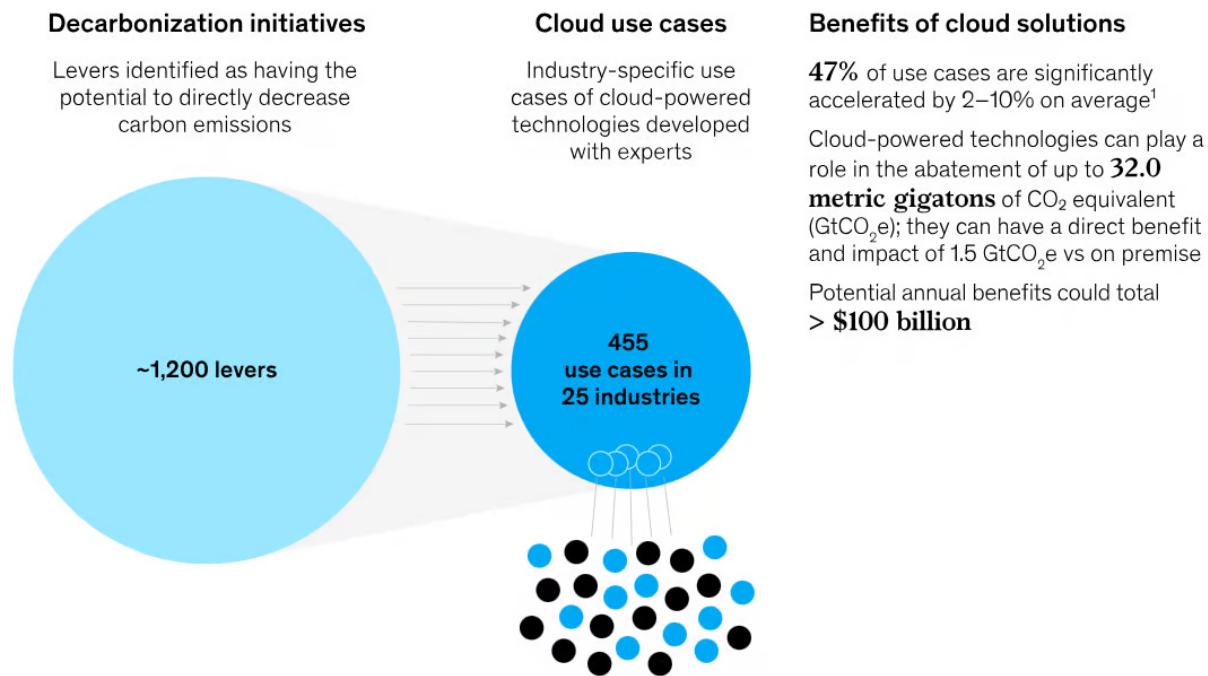
- 1. Stable Energy Supply:** Nuclear energy offers a consistent and reliable power source, essential for the high energy demands of AI technologies.
- 2. Environmental Impact:** By providing a low-carbon energy option, nuclear power helps reduce the carbon footprint associated with AI's intensive energy use.
- 3. Scalability and Economic Efficiency:** Nuclear energy's high energy density and cost-effectiveness make it a scalable solution for the growing energy needs

of AI development.

Some organizations are creating dedicated roles to monitor and reduce AI's environmental footprints by auditing carbon usage across operations. Such efforts highlight how seriously tech companies take these challenges. Integrating sustainability is crucial for AI development.

In response to these challenges, there is a growing focus on sustainable solutions and innovations. One approach is the development of more energy-efficient AI algorithms and hardware. Researchers are exploring ways to reduce the computational complexity of AI models without compromising their performance. This includes techniques like model pruning, quantization, and the use of more efficient neural network architectures.

Cloud-powered technologies will play a central role in the sustainability transition.



Cloud-powered technologies will help with sustainability. Source McKinsey

Industry Initiatives and Corporate Responsibility

Major tech firms investing in AI have sustainability commitments publicly listed. For example, Google aims to solely rely on carbon-free energy by 2030. Microsoft plans to be carbon-negative by 2030 and then remove historical emissions by 2050.

Groups offer guidance on how organizations can implement environmentally responsible AI practices. But more progress across the AI field is required to make systems greener at their core development and infrastructure.

Ten families of climate technologies can play important parts in mitigating carbon emissions.

Climate technology families and examples

				
Renewables Solar, wind (onshore and offshore), grid innovation	Batteries and energy storage Electric-vehicle batteries, long-duration energy storage	Circular economy Battery recycling, chemical cellulosic recycling, heat recovery, plastics recycling	Building technologies Geothermal heating, heat pumps, electric equipment	Industrial-process innovation Electrification of heat sources, green steelmaking, green cement making
				
Hydrogen Electrolyzers, fuel cells, methane pyrolysis	Sustainable fuels Advanced biofuels, e-fuels	Nature-based solutions Monitoring and verification for forests, peatlands, mangroves	Carbon removal, capture, and storage Point-source carbon capture, direct air capture	Agriculture and food Precision agriculture, crop preservation, regenerative tech, alternative proteins

McKinsey & Company

How to mitigate carbon emissions. Source McKinsey

Looking Ahead: The Future of AI and the Environment

Advancing AI must balance with environmental sustainability, requiring efficient computing, renewable energy use, and holistic integration of sustainability in AI development. Innovations in energy-efficient AI, responsible industry practices, and a full understanding of AI's environmental impact are crucial. Prioritizing sustainability ensures AI's positive societal and environmental contributions.

I'm proud to work for IBM that committed to [Net Zero Greenhouse Gas Emissions by 2030](#) with the following plan:

- Reduce its greenhouse gas emissions by 65% by 2025 against the base year 2010. What's most important in the fight against climate change is to actually reduce emissions. The company's net-zero goal is also accompanied by a specific, numerical target for residual emissions that are likely to remain after IBM has first done all it can across its operations to reduce.
- Procure 75% of the electricity it consumes worldwide from renewable sources by 2025, and 90% by 2030.
- Use feasible technologies, such as carbon capture (in or by 2030) to remove emissions in an amount that equals or exceeds the level of IBM's residual emissions.



Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States



Powered by beehiiv

Subject: Day 11: The AI Engineer Profession and Skills
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Saturday, January 13, 2024 2:47:13 PM GMT+05:30
Date Received: Saturday, January 13, 2024 2:47:15 PM GMT+05:30

Day 11: The AI Engineer Profession and Skills

Welcome to Day 11, where we explore the evolving and dynamic role of **AI Engineers** in the rapidly advancing field of Generative AI.

AI Engineering: A New Frontier

- **Role Definition:** AI Engineers are the architects behind practical AI applications, handling everything from the development to the deployment of AI systems.
- **Emerging Importance:** As AI capabilities grow, particularly with the advent of Foundation Models, AI Engineers have transitioned from niche specialists to key players in tech and business landscapes.

I believe the AI Engineer will be the highest-demand engineering job of the decade.

Responsibilities of AI Engineers

- **AI Infrastructure:** Develop and manage robust AI systems, ensuring scalability and efficiency.
- **Advanced Prompting Strategies:** Utilize tools like LangChain for sophisticated prompt engineering with LLMs.
- **Data Management and Model Operations:** Master data preprocessing and embedding techniques, and manage a diverse range of language models for varied applications.
- **AI Model Integration:** Transform AI models into accessible APIs for seamless integration with software systems.

Why AI Engineering is the Future

- **Demand and Recognition:** With AI integration becoming crucial for businesses, AI Engineers are in high demand for their ability to turn AI advancements into practical solutions.
- **Diverse Backgrounds:** Professionals in this field are proving that diverse backgrounds can contribute significantly to AI product development, beyond traditional academic pathways.

The Path to Becoming an AI Engineer

- **Foundational Skills:** Master programming (Python), machine learning, deep learning, and cloud computing.
- **Recommended Courses:** Consider enrolling in courses like IBM AI Engineering Professional Certificate or DeepLearning.ai's specialized courses in AI and LLM application development.
- **Portfolio Development:** Build a portfolio showcasing your AI projects to demonstrate your skills to potential employers.

Here's a good list of courses to start your AI Engineering journey: [link](#)

Armand Ruiz posted this • 2mo



7 Short Courses to Become an AI Engineer in 2024

With 70 days remaining this 2023, you can prepare for next year! ...[show more](#)



53 comments

AI Engineer vs Data Scientist

- **Role Distinction:** Data scientists focus on data analysis and model building, while AI Engineers specialize in building and deploying AI systems and infrastructure.
- **Unique Contribution:** AI Engineers are pivotal in implementing large-scale AI systems, like LLMs, in practical, operational environments.

Factor	AI Engineer	Data Scientist
Responsibilities	Build and deploy AI systems. Work on the underlying infrastructure that powers AI systems.	Collect, clean, analyze, and interpret data. Build and deploy machine learning models.
Skills	Machine learning, programming, cloud computing, distributed systems	Statistics, machine learning, programming, data visualization
Career path	Software engineer → machine learning engineer → AI engineer	Data analyst → data scientist → machine learning engineer → AI engineer
Salary	\$110,000 - \$150,000	\$98,000 - \$137,000
Job outlook	Very good	Very good

Difference between AI Engineer and Data Scientist

The Future of AI Engineering

- **Market Trends:** The demand for AI Engineers is expected to surpass that of data scientists as the field continues to evolve towards more complex, large-scale AI implementations.
- **Shaping the AI Landscape:** AI Engineers are key to crafting strategies and architectures that will define the future of digital innovation, making this role not just current but increasingly vital in the AI-driven future.

AI Engineers stand at the forefront of the AI revolution, turning the promise of AI into tangible, valuable applications. This profession is not just about technical prowess; it's about shaping the future of how we interact with technology.

My prediction:

In numbers, there are probably going to be significantly more AI Engineers than there are data scientists.

Stay tuned for more tomorrow!

Best,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Day 10: The Emergence of Small Language Models
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Friday, January 12, 2024 2:48:35 PM GMT+05:30
Date Received: Friday, January 12, 2024 2:48:37 PM GMT+05:30

Day 10: The Emergence of Small Language Models

On Day 10, we focus on the emerging trend of Small Language Models (SLMs) in the business world and their growing importance alongside Large Language Models (LLMs).

Understanding Large Language Models: A Refresher

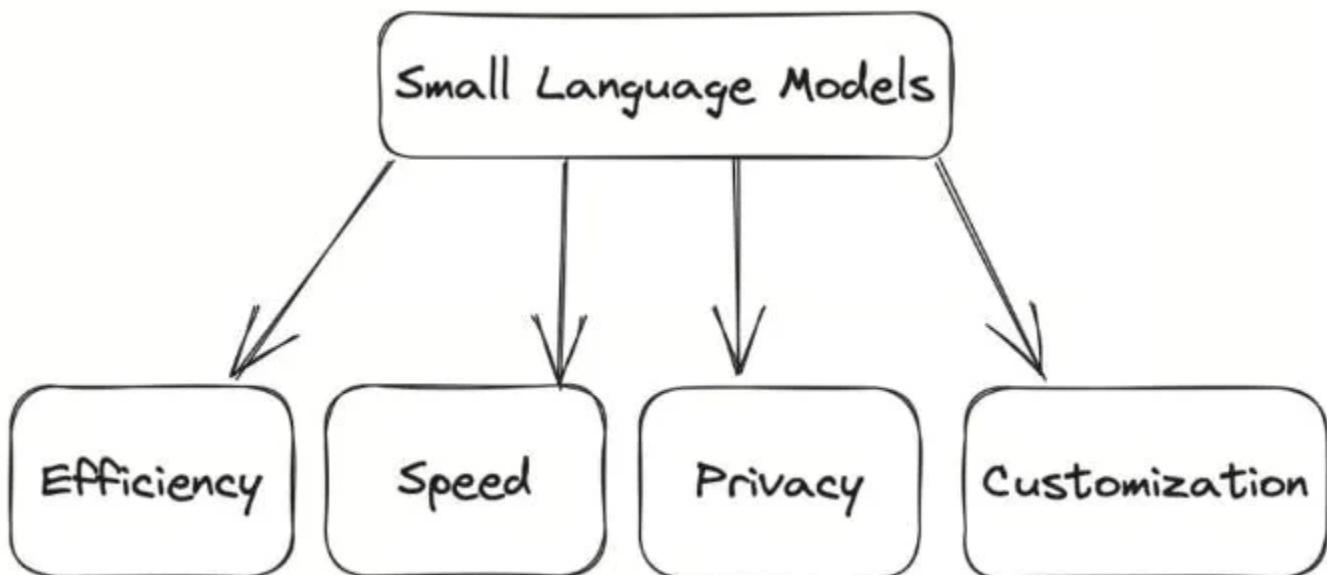
- **Foundation of Today's NLP:** LLMs, trained on vast text data, excel in generating coherent text and performing complex language tasks.
- **Size and Complexity:** Models like GPT-3 (175 billion parameters) and PaLM (540 billion parameters) represent the massive scale of LLMs, offering advanced capabilities but sometimes leading to challenges in accuracy and behavior.

but LLMs are *very* expensive:

running ChatGPT costs approximately **\$700,000 a day**

What are Small Language Models (SLMs)?

- **Defining SLMs:** Generally defined as models with up to 20 billion parameters, SLMs are tailored for specific business tasks like chat, analytics, and content generation.
- **Agility and Customization:** SLMs offer a balance of capability and control, making them well-suited for focused business applications.



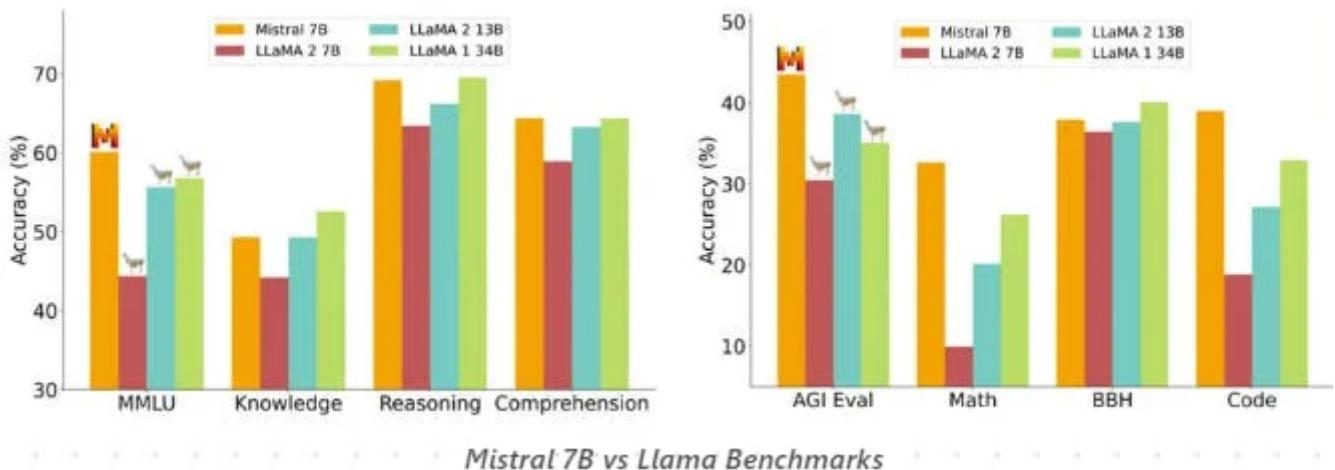
Advantages of SLMs

- **Development and Risk Control:** Easier to build and modify, SLMs reduce risks like bias and hallucinations due to simpler knowledge representations.
- **Efficiency and Sustainability:** Being lightweight and less computationally intensive, SLMs are ideal for deployment on smartphones and edge devices, contributing to sustainability.
- **Cost-Effectiveness:** SLMs offer significant cost savings, making AI more accessible for businesses.

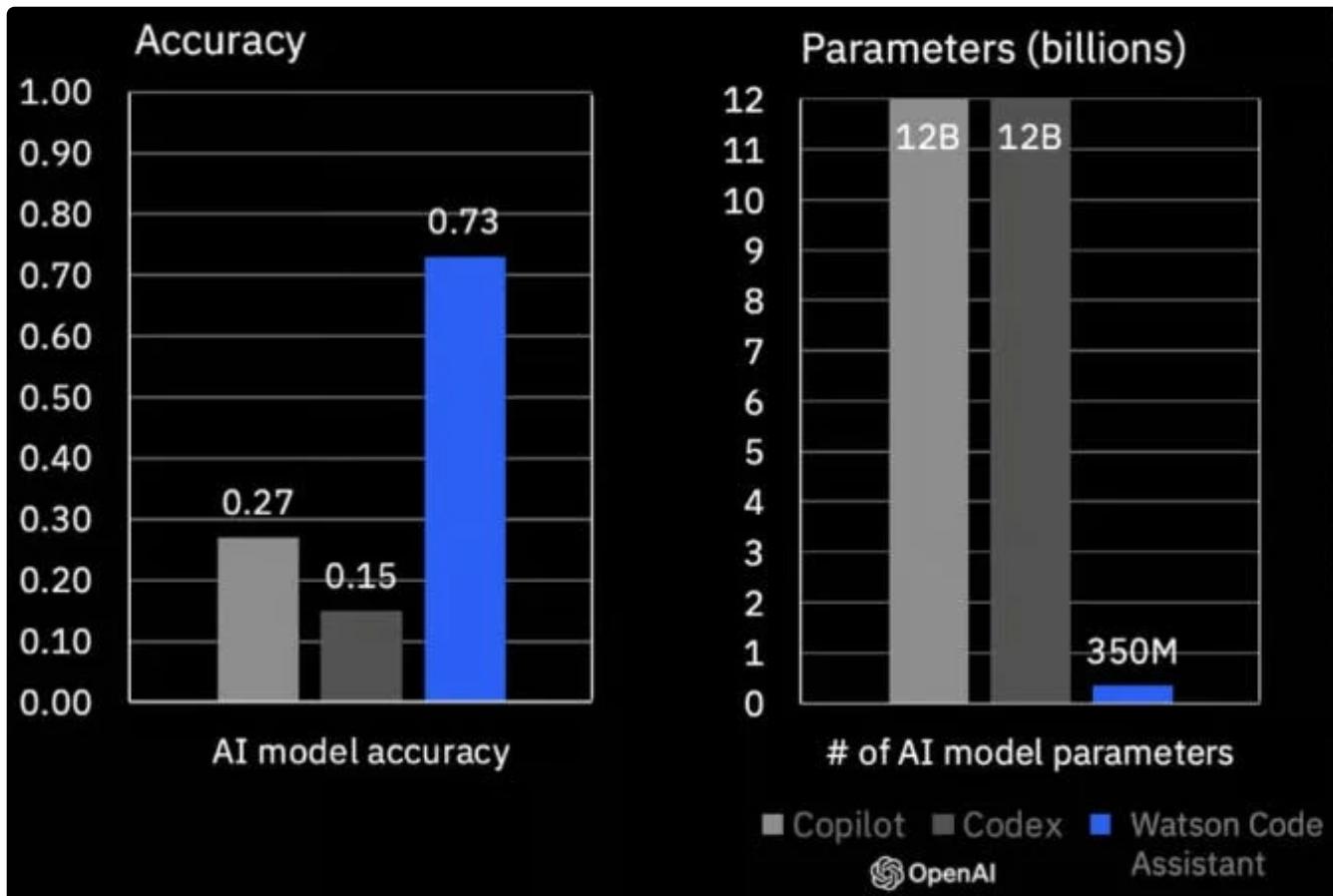
The speed of learning SLMs allow is huge, too. They're within the reach of so many more teams at lower cost. It just lets more innovation cycles happen faster - Brad Edwards

Benchmarking SLMs Against LLMs

- **Performance Comparisons:** For instance, Mistral 7B outperforms larger models in certain benchmarks, demonstrating that SLMs can compete with or even surpass LLMs in specific tasks.



- **Focused Training:** SLMs like IBM Granite, despite smaller size and data, show competitive performance due to targeted training on industry-specific data.



Tuning Small Language Models

- **Customization Techniques:** Similar to LLMs, SLMs can be fine-tuned using various methods to enhance performance for specific use cases.
- **Example of Tuning:** IBM's Granite series, for instance, underwent specialized training for coding, showing how SLMs can be tailored to specific domains.

Use Cases of SLMs

- **Versatile Applications:** SLMs are effective in text generation, chatbots, Q&A, and summarization, offering optimized solutions for resource-limited scenarios.
- **Domain-Specific Tuning:** SLMs can be trained for specialized fields like medical, legal, or technical translation, offering more accurate and relevant outputs than general-purpose LLMs.

In Summary

- **Balancing Capability and Practicality:** SLMs are emerging as a practical alternative to LLMs in many business scenarios, offering a mix of specialized capabilities and control.

SLMs represent a significant development in the AI landscape, providing businesses with more agile, cost-effective, and focused solutions for integrating AI into their operations.

See you tomorrow on Day 12 of this course!

Best,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Day 9: The Generative AI Application Development Stack
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Thursday, January 11, 2024 2:47:27 PM GMT+05:30
Date Received: Thursday, January 11, 2024 2:47:28 PM GMT+05:30

Day 9: The Generative AI Application Development Stack

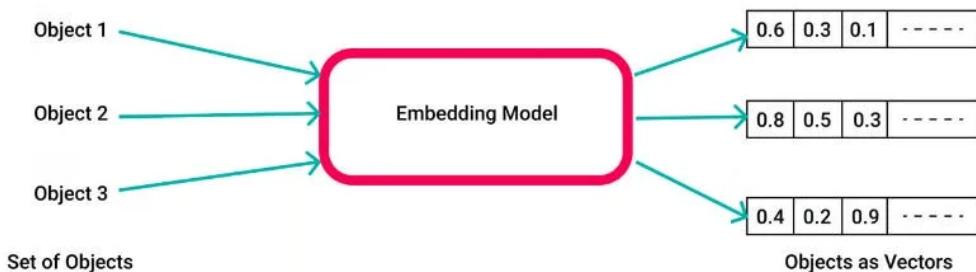
Welcome to Day 9! Today, we're diving into the architecture of the Generative AI (GenAI) stack, crucial for crafting customized GenAI applications.

The GenAI Stack: A Modular, Integrated System

- **Understanding the GenAI Architecture:** This system includes data pipelines, training and inference engines for LLMs, model registries, deployment monitoring, and user interfaces. Tools like LangChain offer orchestration layers for rapid transitions from data to models to apps.

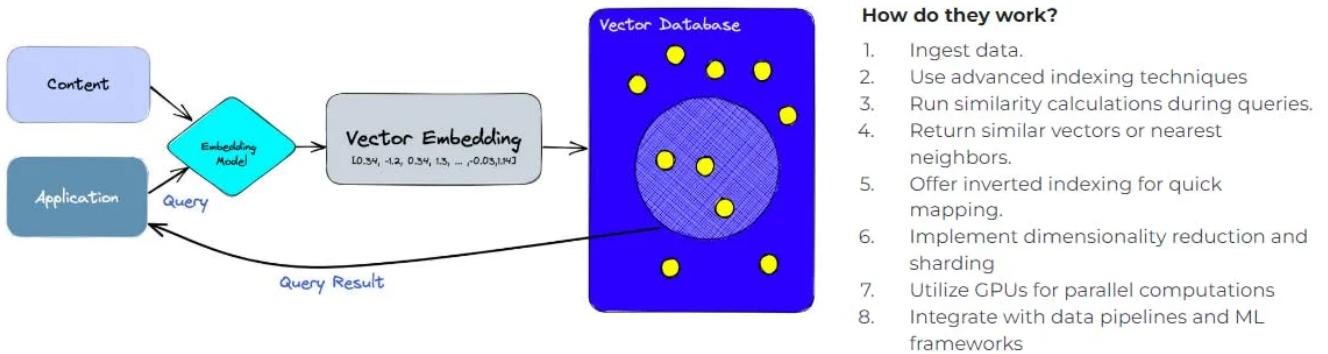
Key Elements of the GenAI Stack:

1- Embeddings (Vectors): These transform high-dimensional data into lower-dimensional vectors, retaining essential information in a more manageable form.

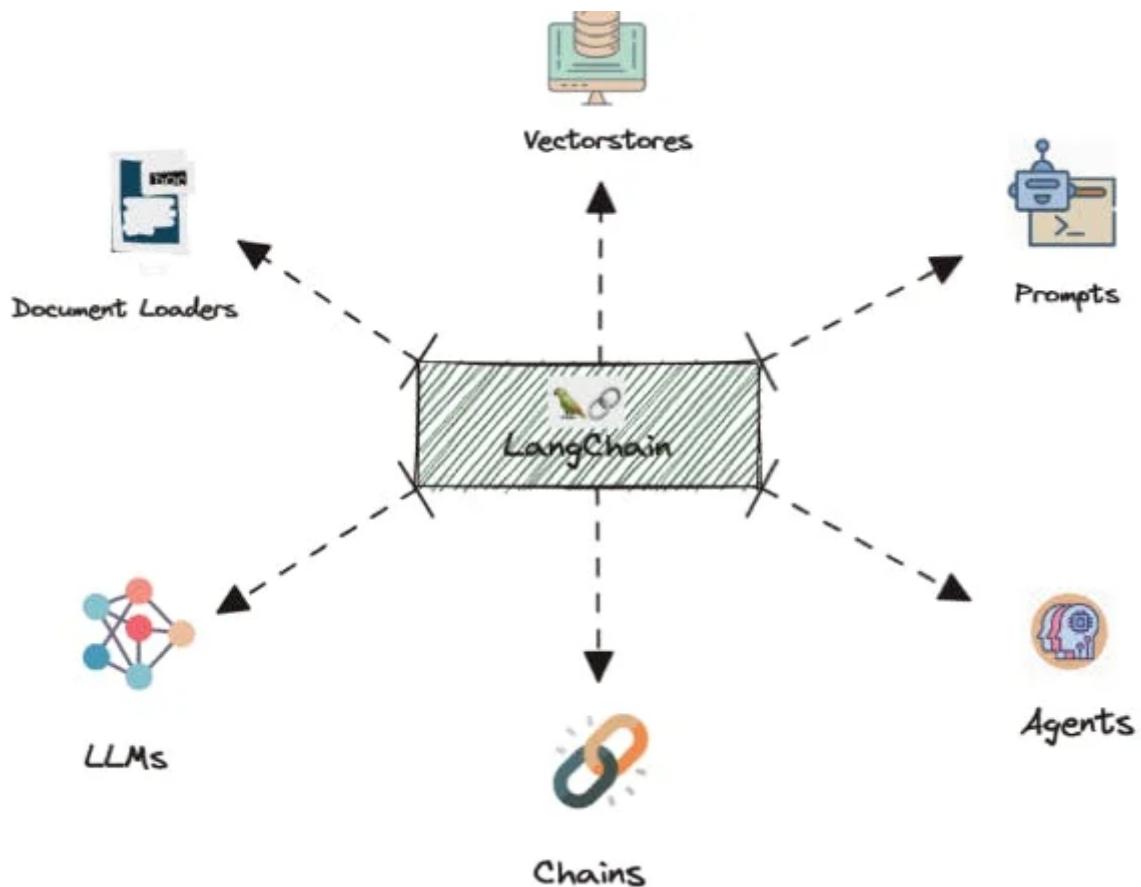


Embeddings convert messy real-world data into mathematical representations capturing hidden relationships. This transformed data powers cutting-edge AI.

2- Vector Database: Stores and indexes vector representations for quick retrieval, supporting operations like vector search and similarity rankings, forming the backbone of vector infrastructure in AI.



3- LangChain: An open-source framework built around LLMs, LangChain facilitates the design and development of various GenAI applications, including chatbots and Generative Question-Answering (GQA).



4- LLMs and Prompts: The core of generative capabilities, LLMs respond to prompts to generate text, making them essential for applications like content creation and customer service.

Building a Simple GenAI App - Step-by-Step:

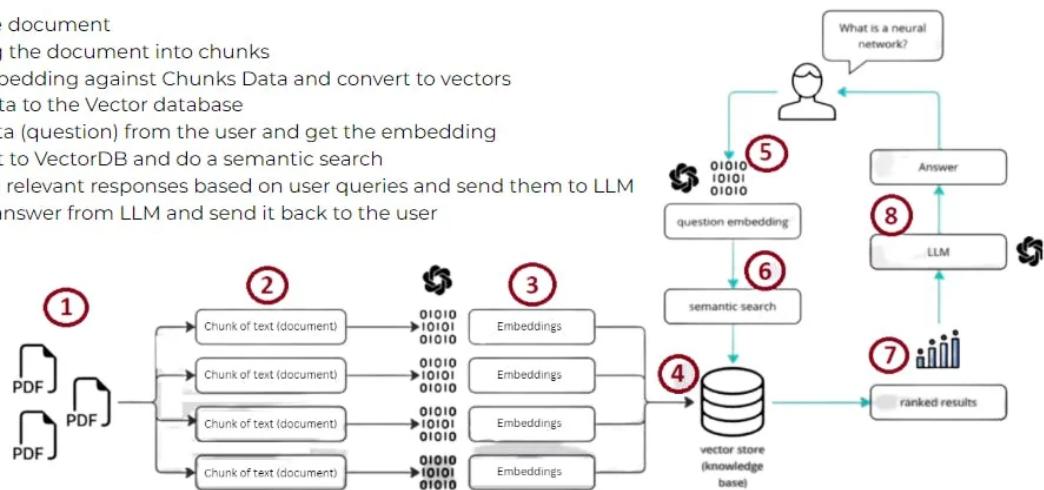
1. **Load Document:** Begin by loading the document or data source.
2. **Split into Chunks:** Break the document into manageable parts.

3. **Create Embeddings:** Convert these chunks into vector representations using embeddings.
4. **Store in Vector Database:** Save these vectors in the database for efficient retrieval.
5. **User Interaction:** Receive queries or input from the user and convert them into embeddings.
6. **Semantic Search in VectorDB:** Connect to the vector database to perform a semantic search based on the user's query.
7. **Retrieve and Process Responses:** Fetch relevant responses, pass them through an LLM, and generate an answer.
8. **Deliver Answer to User:** Present the final output generated by the LLM back to the user.

components of a GenAI solution

Steps:

- Step 1: Load the document
- Step 2: Splitting the document into chunks
- Step 3: Use Embedding against Chunks Data and convert to vectors
- Step 4: Save data to the Vector database
- Step 5: Take data (question) from the user and get the embedding
- Step 6: Connect to VectorDB and do a semantic search
- Step 7: Retrieve relevant responses based on user queries and send them to LLM
- Step 8: Get an answer from LLM and send it back to the user



Understanding and utilizing the components of the GenAI stack is key for businesses looking to leverage AI for innovative applications. This modular approach allows for customization and scalability, fitting various business needs and goals.

Tomorrow, we will talk about Small Language Models!

Cheers,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Day 8: Generative AI Applications and Use Cases
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Wednesday, January 10, 2024 2:48:30 PM GMT+05:30
Date Received: Wednesday, January 10, 2024 2:48:32 PM GMT+05:30

Day 8 - Generative AI Applications and Use Cases

On Day 8, let's talk about the transformative applications of generative AI in business, examining how these technologies are reshaping various industries and enterprise functions.

Broad Spectrum of Business Applications:

- **Marketing and Advertising:** From crafting compelling ad copy to generating creative landing pages, generative AI is revolutionizing how businesses approach marketing.
- **Content Creation:** AI is now capable of producing news articles and social media content, enhancing digital presence with minimal effort.
- **Customer Service:** By deploying AI-driven chatbots, businesses can ensure engaging, natural conversations with customers, elevating the service experience.
- **Summarization:** Generative AI can distill lengthy reports and papers into concise, informative summaries, aiding decision-making and research.
- **Data Analysis:** AI tools can sift through vast datasets, uncovering patterns and insights that drive strategic decisions.
- **Personalization:** Tailoring content to individual user preferences or customer segments is now more efficient with AI's generative capabilities.
- **Product Development:** Rapid prototyping and testing of new product designs are made possible, speeding up the innovation process.

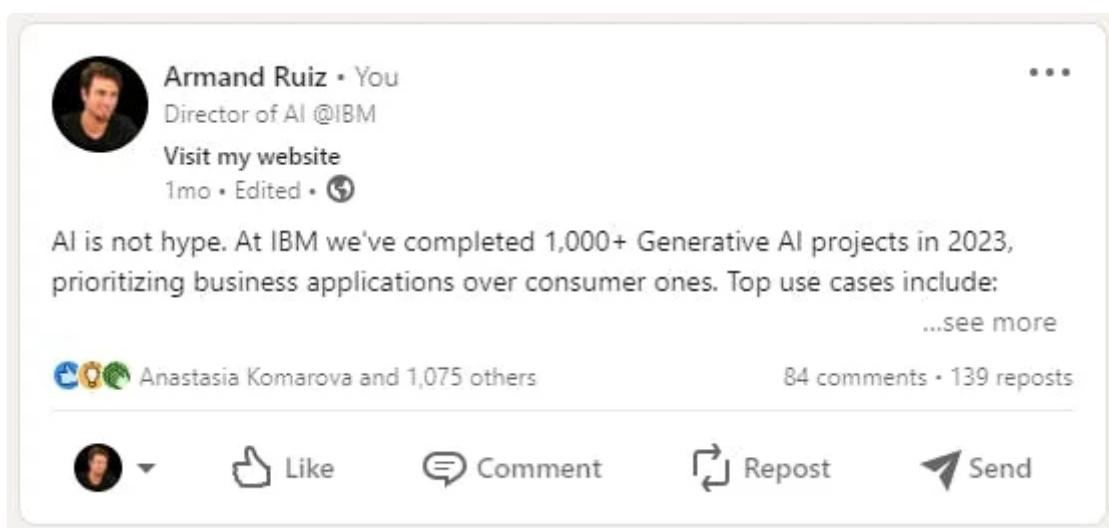
most common generative AI tasks for business

Retrieval-Augmented Generation Based on a document or dynamic content, create a chatbot or question-answering feature. <i>Building a Q&A resource from a broad knowledge base, providing customer service assistance.</i>	Summarization Condenses textual information into concise summaries, like summarizing customer feedback. <i>Conversation summaries, insurance coverage, meeting transcripts, contact information.</i>	Content Generation Generate text content for a specific purpose. <i>Marketing campaigns, job descriptions, blog posts and articles, email drafting support.</i>
Named Entity Recognition Identify and extract essential information from unstructured text. <i>Audit acceleration, SEC 10k fact extraction</i>	Insight Extraction Analyze existing unstructured text content to surface insights in specialized domain areas. <i>Medical diagnosis support, user research findings.</i>	Classification Read and classify written input with as few as zero examples. <i>Sorting of customer complaints, thread and vulnerability classification, sentiment analysis.</i>

Key Benefits for Businesses:

- **Increased Efficiency:** Automate repetitive content generation tasks, freeing up valuable time for strategic work.
- **Cost Savings:** Reduce reliance on expensive human labor, particularly in creative and writing tasks.
- **Consistency and Quality:** Maintain a consistent brand voice while leveraging AI's analytical capabilities to produce high-quality content.
- **Faster Ideation and Scalability:** Accelerate the process of brainstorming and content production, enabling businesses to reach larger audiences more effectively.

See my detailed list of Use Cases in this recent LinkedIn post. Here's the [link](#)



Armand Ruiz • You
Director of AI @IBM
[Visit my website](#)
1mo • Edited •

AI is not hype. At IBM we've completed 1,000+ Generative AI projects in 2023, prioritizing business applications over consumer ones. Top use cases include:
[...see more](#)

 Anastasia Komarova and 1,075 others 84 comments • 139 reposts

  Like  Comment  Repost  Send

Transforming Enterprise Operations: Generative AI is not just a tool; it's a **game-changer** for business operations. It enables businesses to scale creativity, streamline content creation, engage customers in novel ways, and achieve significant time and cost savings. The integration of AI into these functions is transforming how businesses interact with their customers, manage their internal processes, and innovate in their product offerings.

Incorporating generative AI into business strategies can lead to more efficient, creative, and data-driven approaches, opening new avenues for growth and competitive advantage.

That's all for today. Get ready for more insights tomorrow!

Cheers,

Armand ❤

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768, Seattle, WA 98104-2205

Subject: Day 7: The most popular LLMs available
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Tuesday, January 9, 2024 2:47:30 PM GMT+05:30
Date Received: Tuesday, January 9, 2024 2:47:31 PM GMT+05:30

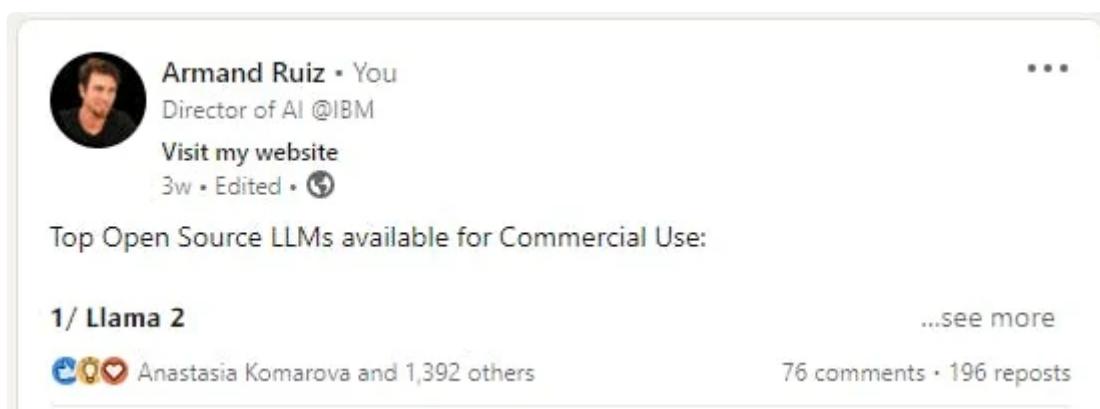
Day 7 - The most popular LLMs available

On Day 7 of our AI series, let's navigate the world of Large Language Models (LLMs) by understanding the key differences between open-source and proprietary models and services, and exploring some of the most popular LLMs available today.

Open Source LLMs:

- **Accessible and Collaborative:** These models are freely available for use, modification, and distribution, promoting community-driven development and innovation.
- **Examples of Open Source LLMs:**
 - **GPT-Neo/GPT-J:** Developed by EleutherAI, these models are open-source alternatives to OpenAI's GPT models, offering similar capabilities.
 - **BERT:** Developed by Google, BERT has been a groundbreaking model for understanding context in natural language, widely used in various applications.

I listed the Top Open Source LLMs a few weeks back on LinkedIn. Here's the [link](#)



Armand Ruiz • You
Director of AI @IBM
[Visit my website](#)
3w • Edited • 

Top Open Source LLMs available for Commercial Use:

1/ Llama 2 ...see more

 Anastasia Komarova and 1,392 others 76 comments • 196 reposts

Closed Source LLMs:

- **Commercial and Proprietary:** These models are developed and maintained by private entities, often requiring licenses or subscriptions for access.
- **Examples of Closed Source LLMs:**

- **OpenAI's GPT-3/GPT-4:** Known for their advanced capabilities, these models have set benchmarks in generative AI but are accessible mostly through API with usage costs.
- **Google's LaMDA:** A cutting-edge model designed for conversational AI, used internally by Google.

Open Source vs Closed Source

Open Source LLMs:

- Language models with publicly available source code.
- Can be freely accessed, used, modified, and distributed by anyone.
- Promote collaboration, transparency, and community involvement.
- Allow for collective development, innovation, and knowledge sharing.

Closed Source LLMs:

- Source code is not publicly available.
- Developed and maintained by private organizations or companies.
- Often commercial products requiring licenses or subscriptions.
- Architecture, training data, and algorithms are typically proprietary and not disclosed to the public.

The Game Changer: Llama 2

- **Accessibility and Versatility:** Meta's Llama 2 has been released as an open-source AI model, making it accessible for everyone from startups to researchers. Its availability in different sizes (7B, 13B, 70B-parameter models) offers a range of options for fine-tuning and deployment.
- **Innovation and Privacy:** As an open-source model, Llama 2 removes barriers to AI adoption and addresses data privacy concerns by allowing private hosting and customization with your own data.
- **Performance Benchmarking:** Llama 2 stands on par with models like GPT-3.5 in terms of performance, particularly excelling in generating helpful responses for prompts. However, it shows less proficiency in coding tasks compared to other specialized models.
- **Cost and Community Benefits:** Meta's open-sourcing of Llama 2, despite the substantial development cost, taps into the collective wisdom of the global AI community, accelerating innovation and potentially leveling the playing field against closed-source counterparts.

Why the Distinction Matters: Understanding the differences between open source and closed source LLMs is crucial for businesses and developers. Open source models offer transparency and the opportunity for customization, while closed source models, often

backed by significant resources and research, provide robust, state-of-the-art capabilities but with usage restrictions and costs.

In your AI endeavors, choosing between open source and closed source LLMs will depend on your specific needs, resources, and goals.

Stay tuned for more exciting developments in AI!

See you tomorrow

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

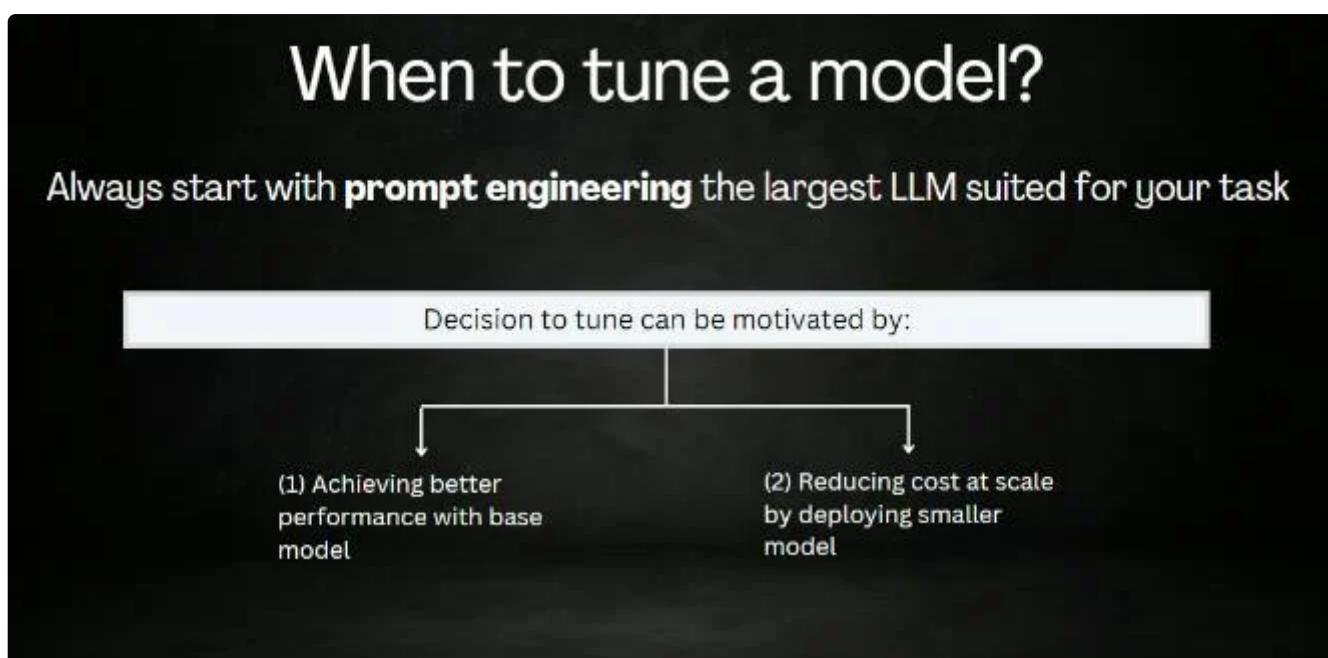
Subject: Day 6: How to customize foundation models
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Monday, January 8, 2024 2:48:12 PM GMT+05:30
Date Received: Monday, January 8, 2024 2:48:13 PM GMT+05:30

Day 6 - How to customize foundation models

Today, on Day 6 of our AI journey, let's unlock the secrets of customizing foundation models to suit your specific needs. Understanding when and how to tune these models is crucial for optimal performance.

Deciding When to Tune Your Model

Starting Point: Begin with prompt engineering using the largest suitable Language Model (LLM) for your task to gauge if LLMs can handle it. Experiment with various prompt formats and examples.



Prompting Techniques:

1. Zero-Shot Prompting:

- **Efficiency with No Extra Data:** This involves giving a natural language prompt to generate desired outputs without additional training data.
- **Example:** "Provide a summary of the following passage: [insert text]."

2. One-Shot Prompting:

- **A Single Example to Guide:** Introduce one example along with your prompt to demonstrate the desired outcome.
- **Example:** "Write marketing copy for WorkoutFuel protein shakes in an enthusiastic, punchy voice," along with a high-energy example text.

3. Few-Shot Prompting:

- **Leveraging a Few Examples:** Provide a handful of examples to establish the pattern or style for the model to replicate.
- **Example:** To generate meeting summaries, give 2-3 examples before asking the model to create new ones.

Technique	Advantages	When to Use
Zero-Shot Prompting	- Simplest to implement - No training data needed	- Quickly test model capabilities - Simple tasks and inferences
One-Shot Prompting	- Can teach complex behaviors - Minimal training data	- Teaching nuanced or contextual responses - When only one good example is available
Few-Shot Prompting	- Teach more robust behaviors - Still very sample efficient	- Targeted enhancement for specific skills - When limited data for a task
Fine-Tuning	- Maximize model performance - Learn complex and nuanced behaviors	- Specializing model for dedicated tasks - When abundant training data exists
Parameter-Efficient Fine-Tuning	- Improved generalization - Faster and lower resource fine-tuning	- Specializing with limited data - Personalizing models for clients

Data-Driven Tuning for Deeper Customization

- **Fine-Tuning:** Adjusting model weights on a specific dataset to cater to your unique objectives, like customizing tone or addressing complex prompts.
- **Parameter-Efficient Fine-Tuning (PEFT):** Delta tuning updates only a small subset of parameters, offering a faster, cost-effective alternative to traditional fine-tuning.

Fine-tuning vs PEFT

Fine-tuning

Tune **ALL** model parameters

Generate a copy of the base model that **requires hosting**

Requires **1,000s - 100,000s** labeled data points

Significant performance gains on target task compared to base model

Prone to catastrophic forgetting

Parameter-efficient fine-tuning (PEFT)

Tune a **small number** of (extra) model parameters

Generates **tiny checkpoints** worth a few MBs or less

Requires **100s - 1,000s** labeled data points

Comparable to full fine-tuning depending on base model size and data used

Overcomes catastrophic forgetting

PEFT Techniques:

Parameter-Efficient Fine-Tuning (PEFT) is a more cost-effective and efficient method because it focuses on optimizing a small subset of model parameters, reducing computational resources and training time while maintaining high-performance levels. There are multiple techniques:

- **Prefix Tuning:** Attaches vectors with free parameters to input embeddings, training them while keeping the LLM frozen.
- **Prompt Tuning:** A simpler variant of prefix tuning, adding a vector only at the input layer.
- **P-Tuning:** Automates the search and optimization of prompts using an LSTM model.
- **LoRA:** Low-Rank Adaptation adds update matrices to existing weights, training these new weights.

Choosing the Right Technique:

- **Goal-Oriented Approach:** Select the customization method based on your specific goals and the data you have. For instance, zero-shot and few-shot prompting work well with minimal data, while data-driven tuning is ideal for more complex, data-rich tasks.

Customizing foundation models can significantly enhance their performance on specific tasks, making them more aligned with your business objectives.

Stay tuned for more insights tomorrow!

Best,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Day 5: What it Takes to Train a Foundation Model
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Sunday, January 7, 2024 2:47:10 PM GMT+05:30
Date Received: Sunday, January 7, 2024 2:47:12 PM GMT+05:30

Day 5 - What it Takes to Train a Foundation Model

Welcome to Day 5 of our AI journey! Today, I will focus on the reasons why training your foundation model can be a pivotal step for your business but it is a decision that has to be taken very wisely.

Control Over the Model

- **Tailored Approach:** When you train your model, you control the data and parameters, allowing you to tailor it to specific styles or domains. This customization ensures the model aligns perfectly with your business needs.

Improved Performance

- **State-of-the-Art Results:** Foundation models trained on large, diverse datasets can outperform pre-trained models, especially if your dataset is domain-specific.

Customization

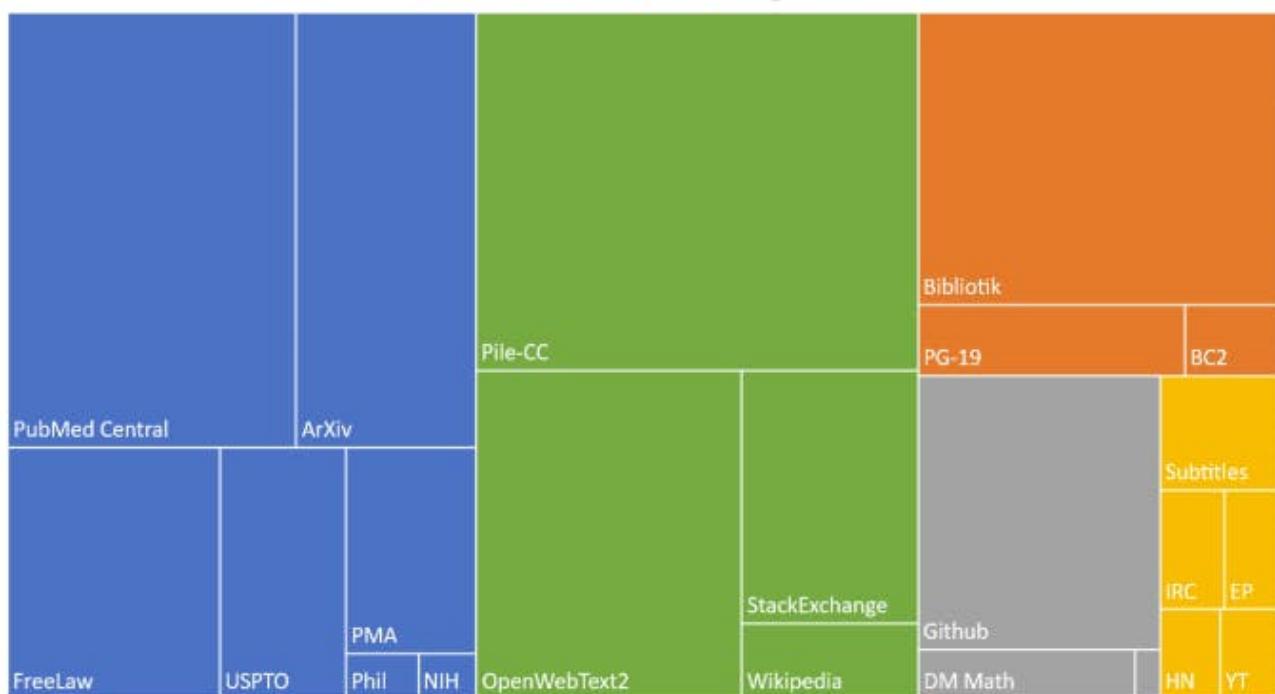
- **Modifying Architecture:** Building your own model means you can alter aspects like tokenizer, vocabulary size, or model architecture, which might be necessary if these components are central to your business strategy.

Challenges of Training from Scratch

- **Data Collection:** Amassing a large, relevant dataset is crucial. An example is The Pile, an extensive, diverse language modeling dataset.

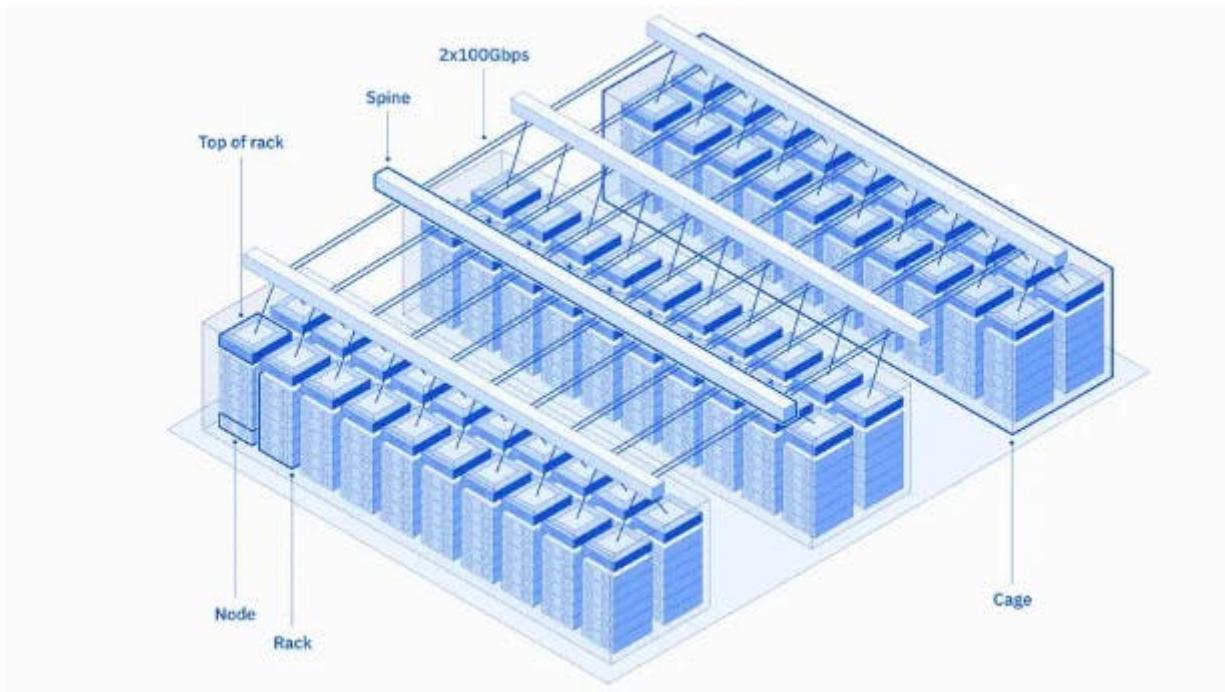
Composition of the Pile by Category

■ Academic ■ Internet ■ Prose ■ Dialogue ■ Misc



The PILE, public data used to train LLMs

- **Compute Resources:** Significant computational power is needed, as demonstrated by AI supercomputers, equipped with thousands of Nvidia A100 and H100 GPUs.



IBM AI Supercomputer VELA

- **Expertise:** Specialized knowledge in AI and ML is essential due to the complexity of model architecture and training processes.

 OpenAI is setting new industry standards with its engineers earning an average annual salary of \$925,000! This includes a base salary of \$300,000 and a whopping \$625,000 in stock-based compensation. Some even earn as much as \$1.4 million!  #AI #OpenAI #TechNews

Training Steps

1. **Dataset Collection:** Gather a large, diverse dataset relevant to your tasks.
2. **Preparation and Tokenization:** Clean and format your data, breaking down text into tokens.
3. **Configure Training:** Set hyperparameters, choose the architecture, and allocate computational resources.
4. **Training:** Train your model using deep learning algorithms.
5. **Evaluation:** Test the model's performance on a separate dataset.
6. **Deployment:** Once satisfied, deploy the model for practical use.

Cost Considerations

Training a foundation model can range from tens of thousands to millions of dollars, depending on the model size, data volume, and computational resources.

Recommendation for Businesses

- **Customize a Pre-Trained Model:** Starting with a pre-trained model and customizing it with techniques like Parameter Efficient Fine Tuning (PEFT) can save time and resources.
- **Consider Needs and Resources:** Evaluate your specific needs and available resources to decide between purchasing, training, or customizing a model.

Customizing foundation models is a great way to get the most out of these powerful tools. It is less expensive, faster, and can give you better performance than training a model from scratch.

In conclusion, while training a foundation model is resource-intensive, it offers unparalleled control and performance, essential for businesses aiming to develop a strong technological edge in AI.

Tomorrow, we'll explore another exciting aspect of AI.

Best,

Armand ☺

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Make 2024 the year of Learning AI
From: "Armand from nocode.ai" <nocodeai@mail.beehiiv.com>
To: "deepakchawla35@gmail.com" <deepakchawla35@gmail.com>
Date Sent: Saturday, January 6, 2024 5:31:36 PM GMT+05:30
Date Received: Saturday, January 6, 2024 5:31:37 PM GMT+05:30

January 06, 2024 | [Read Online](#)

Make 2024 the year of Learning AI

no matter your profession



Welcome to the 3,654 new members this week! **nocode.ai** now has 28,358 subscribers



Happy 2024, folks! Let's make this year the one where we all learn the basic concepts of AI and put them into practice in our daily lives and jobs.

In this first edition of the year, I'm sharing with you **FREE** assets to make 2024 the year you break into the magical world of AI, regardless of your previous experience.

Today, I'll cover:

- My NEW Generative AI Email Course
- The AI Bootcamp
- Top AI Ted Talks
- Coursera Generative AI Specialization from IBM

Let's dive in!



1. My New Generative AI Email Course

LLMs, SLMs, NLP, Embeddings, Vector Databases, Tokens, LLMOps, LVMs, RAG,

Are you feeling lost trying to understand all these AI terms?

I've created a 15-day course to help you know all the core concepts of Generative AI, in just 5 minutes a day.

These are the topics I cover:

- **Day 1: Introduction to Generative AI** Overview of generative AI and its importance in business.
- **Day 2: Types of Generative AI Models** Exploring different generative AI models like GANs, VAEs, and transformers.
- **Day 3: Traditional ML vs Generative AI** Comparing traditional machine learning with generative AI methods.
- **Day 4: What are GPUs** Understanding the role of GPUs in AI and machine learning tasks.
- **Day 5: What it Takes to Train a Foundation Model** Insights into the resources and processes for training large foundation models.
- **Day 6: How to Customize Foundation Models** Discussing techniques for customizing foundation models for specific uses.
- **Day 7: The Most Popular LLMs Available** - Overview of the most widely-used large language models and their features.
- **Day 8: Generative AI Applications and Use Cases** Exploring practical applications of generative AI across business sectors.
- **Day 9: The Generative AI Stack** Understanding the components and architecture of the generative AI tech stack.
- **Day 10: The Emergence of Small Language Models** Discussing the rise and importance of small language models in AI.
- **Day 11: The AI Engineer Profession and Skills** Exploring the role, responsibilities, and required skills of AI engineers.
- **Day 12: Ethical Considerations in AI** Discussing the ethical challenges in AI development and deployment.
- **Day 13: Create Your AI for Business Roadmap** How to develop a strategic AI integration roadmap for businesses.
- **Day 14: Future Trends in AI** Exploring future developments and trends in AI.

• Day 15: Continuing Your AI Journey

How does it work?

1. Enroll here: <https://go.nocode.ai/30-day-email-course>
2. Receive a short lesson via email every day for the next 15 days after enrollment.
The content is designed to be completed in 3-5 minutes.
3. Discuss what you learn with your colleagues and peers, and consider how the new concepts relate to your business and current job.

This is your best opportunity to get ahead and gain essential skills. Get 1% better every day. Knowledge compounds over time.

Learn GenAI in 15 days

2. The AI Bootcamp

In case you want to invest more time the AI Bootcamp will take you From Zero to Hero, Learn the Fundamentals of AI. This Free Course is perfect for beginner to intermediate-level professionals who want to break into AI. Transform your skillset and accelerate your career.

This course includes:

- **50 videos:** From the basics to advanced concepts and real demos. I explain everything step-by-step.
- **15 Practice Exercises:** Learn by doing. Train and deploy your first models to understand the core concepts. No technical experience is required.
- **Learn in Community:** Be part of the private group where I will be answering your questions and providing my top tips.

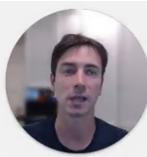
It is ideal for:

- ✓ Professionals looking to break into AI and transform business
- ✓ Learners new to AI seeking to successfully experiment
- ✓ Anyone eager to grow their AI network and connect with experts

I created a very complete course curriculum:

- Module 1: AI Basics and Core Concepts
- Module 2 - It's all about the Data
- Module 3 - Traditional Machine Learning
- Module 4 - The Deep Learning Revolution
- Module 5 - Put AI into Production
- Module 6 - Generative AI
- Module 7 - Ethical and Trustworthy AI
- Module 8 - AI Tools
- Module 9 - Implement AI for Business
- Module 10 - A Sneak Peek into the Future

It includes step-by-step labs to learn how to prepare data, image classification, text classification, fine-tune LLMs, and develop chatbots, AI Agents, and AI Avatars. A lot of fun!



the labs

exercise 1 brainstorm ai use cases for your business <small>with Clarifai</small>	exercise 2 Data Transformation and Automation <small>with Tabula</small>	exercise 3 Chat and Explore your Data <small>with Akkio</small>	exercise 4 Train your first ML Model <small>with IBM AutoAI</small>	exercise 5 Image classification <small>with Clarifai</small>
exercise 6 text classification <small>with Clarifai</small>	exercise 7 Put your first model into production <small>with IBM Watson Studio</small>	exercise 8 Getting started with GPT-4 <small>with OpenAI</small>	exercise 9 Create your first chatbot <small>with Stack-AI</small>	exercise 10 Chat with your data and documents <small>with Stack-AI</small>
exercise 11 fine-tune your first LLM <small>with MontsterAPI</small>	exercise 12 build a trustworthy ai system <small>with IBM OpenScale</small>	exercise 13 share ai projects gone wrong	exercise 14 generate images & create your avatars <small>with StableDiffusion</small>	exercise 15 Create your first AI Agent <small>with LangChain</small>

all using no-code ai tools!

the ai bootcamp

[Join the AI Bootcamp](#)

3. Top AI Ted Talks

10 AI TED Talks. Watch them all and get inspired for the new year!

List of Ted Talks —> [link](#)

Armand Ruiz • You
Director of AI @IBM
[Visit my website](#)
2w •

10 AI TED Talks for the remaining days of 2023. Watch them all and get inspired for the new year!

1. The exciting, perilous journey toward AGI

Before a management change at OpenAI, cofounder Ilya Sutskever examined AGI's transformative potential.

<https://bit.ly/4aloSL4>

4. Coursera Generative AI Specialization from IBM

IBM Coursera. IBM just released the Generative AI Fundamentals Specialization. You can enroll now and this is what you will learn:

[Learn more](#)

Happy learning everyone. I wish you a wonderful 2024!

X in

Update your email preferences or unsubscribe [here](#)

© 2024 the nocode.ai newsletter

228 Park Ave S, #29976, New York, New York 10003, United States

Powered by beehiiv

Subject: Day 4: Introduction to GPUs
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Saturday, January 6, 2024 2:47:20 PM GMT+05:30
Date Received: Saturday, January 6, 2024 2:47:22 PM GMT+05:30

Day 4 - Demystifying GPUs in AI

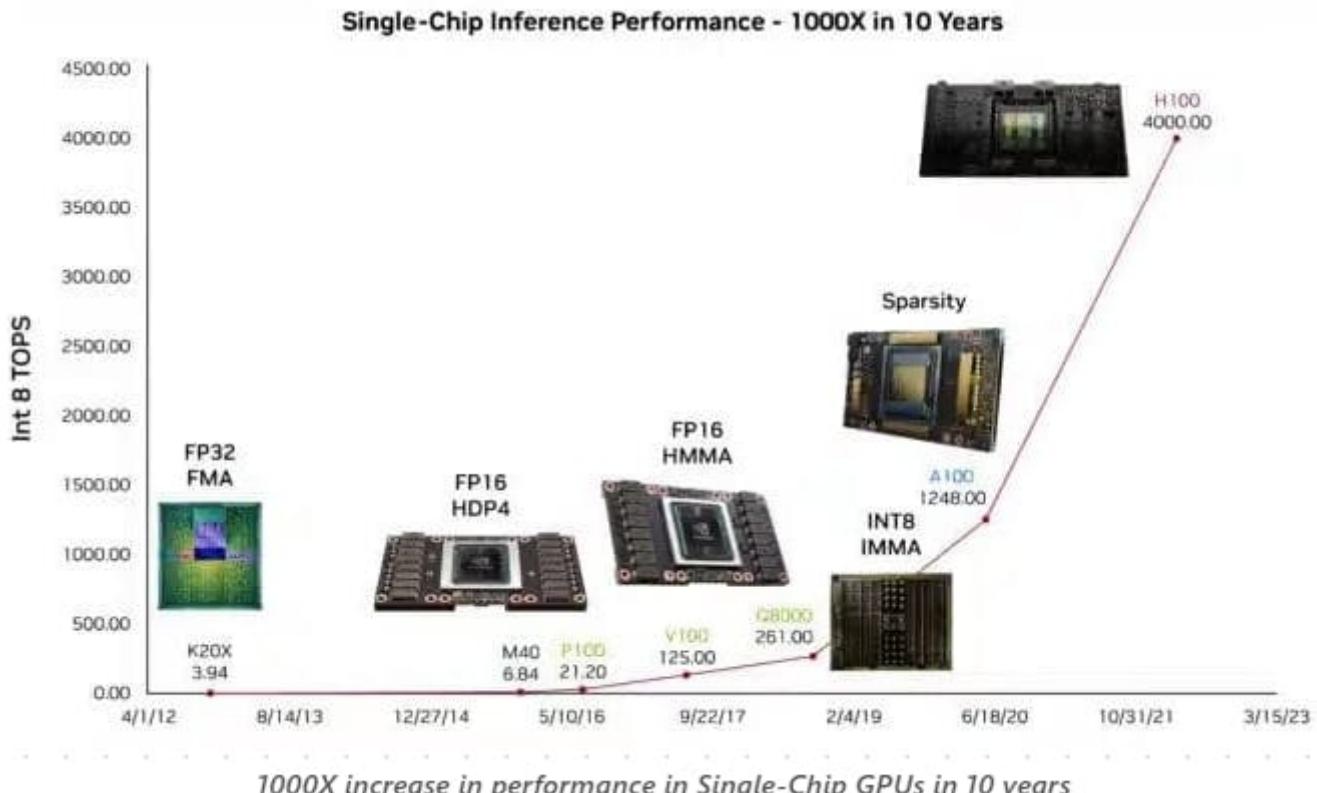
Welcome to Day 4! Today, we're focusing on **GPUs** - the driving force behind the AI revolution, especially in training and inference tasks.

What are GPUs?

GPUs, or Graphics Processing Units, were initially designed for rendering graphics and video tasks. Remember the excitement of upgrading your PC with a new GPU for better gaming? That same technology has become pivotal in AI, thanks to its architecture and parallel processing capabilities.

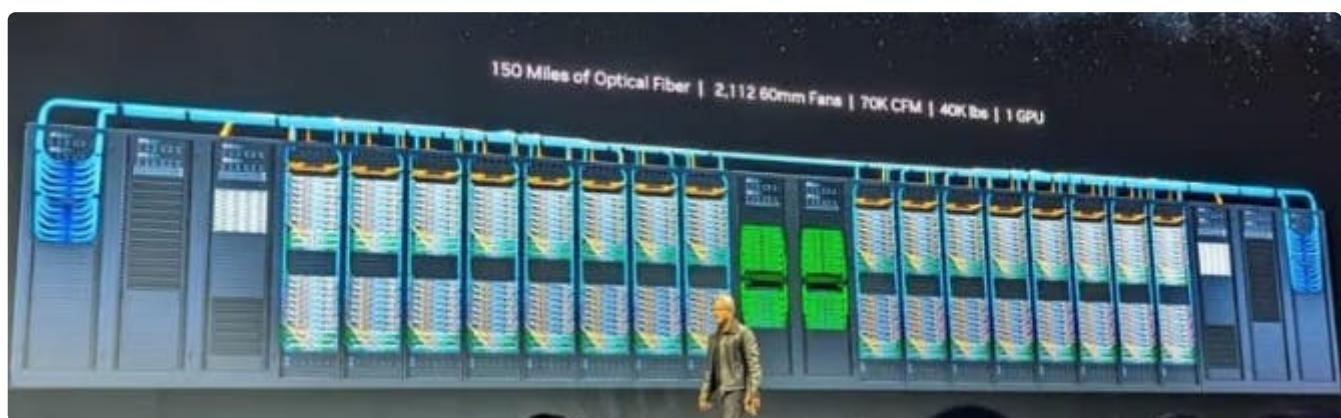
Why are GPUs Crucial for AI?

- **GPU Architecture:** GPUs have thousands of small cores (organized into Streaming Multiprocessors) designed for parallel computing. This setup is perfect for AI workloads, which often require simultaneous processing of large data sets.
- **Parallel Processing Power:** GPUs excel at performing multiple calculations at once, making them ideal for handling the complex mathematical operations needed in AI.
- **Speed and Efficiency:** With rapid thread switching and high memory latency tolerance, GPUs significantly reduce the time needed for training neural networks, like GPT-3, which requires 300 zettaflops of computing power.
- **AI Framework Support:** Manufacturers like Nvidia, AMD, and Intel have optimized GPUs for AI frameworks such as TensorFlow and PyTorch.



The Rise of AI Supercomputers

AI research has skyrocketed with the advent of AI supercomputers, clusters of GPUs working together. These supercomputers, like Summit, Sierra, and Fugaku, are pushing the boundaries in fields like scientific research and climate modeling.



CEO of Nvidia presenting DGX GH200 AI Supercomputer

Selecting the Right GPU

Choosing the appropriate GPU involves understanding your needs - whether it's for AI training, inference, or both, and balancing factors like budget, performance, compatibility, and scalability.

The Future of GPUs in AI

We're witnessing a 1000X increase in single-chip GPU performance over the last decade. The future holds more specialized AI chips, quantum computing advancements, and edge AI integrations, further transforming the AI landscape.

In Conclusion, GPUs have transitioned from enhancing our gaming experiences to becoming the backbone of AI, fueling advancements in machine learning and beyond. Their role in accelerating AI workloads is indisputable and will continue to shape the future of technology.

Tomorrow, we'll explore another exciting aspect of AI. Stay tuned!

Best,

Armand

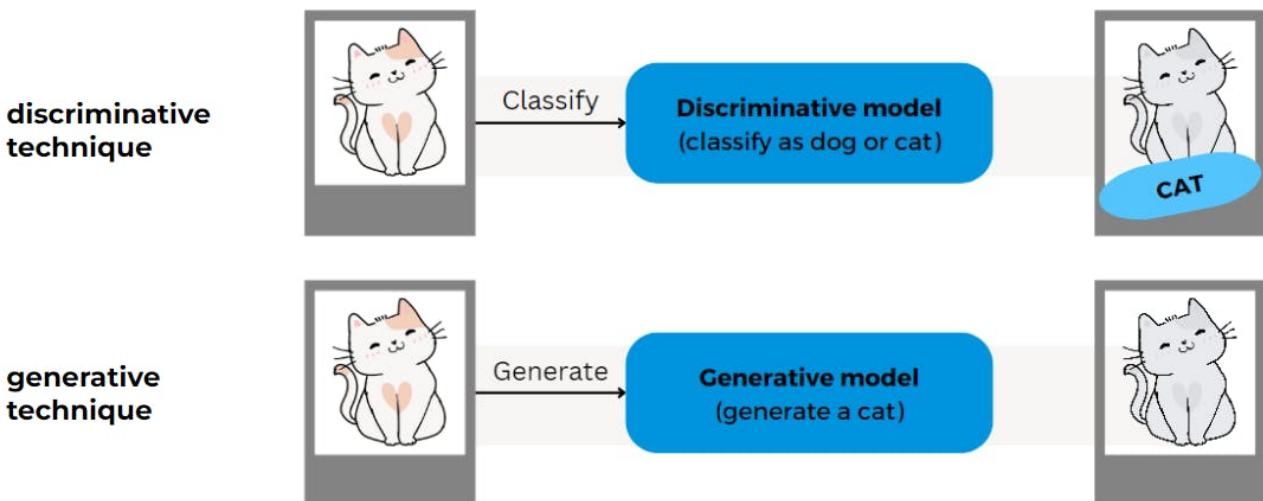
[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Day 3: Traditional ML vs Generative AI
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Friday, January 5, 2024 2:47:31 PM GMT+05:30
Date Received: Friday, January 5, 2024 2:49:03 PM GMT+05:30

Day 3 - Traditional ML vs Generative AI or Discriminative vs Generative Models

Welcome to Day 3! Today, we're diving deeper into the AI landscape by contrasting Traditional Machine Learning (ML), focusing on discriminative models, against Generative AI, which revolves around generative models.

To start, let's understand the difference between with a simple example:



Traditional Machine Learning

Discriminative models in Traditional ML are designed to classify or predict outcomes based on input data. They focus on drawing boundaries between different categories and making decisions.

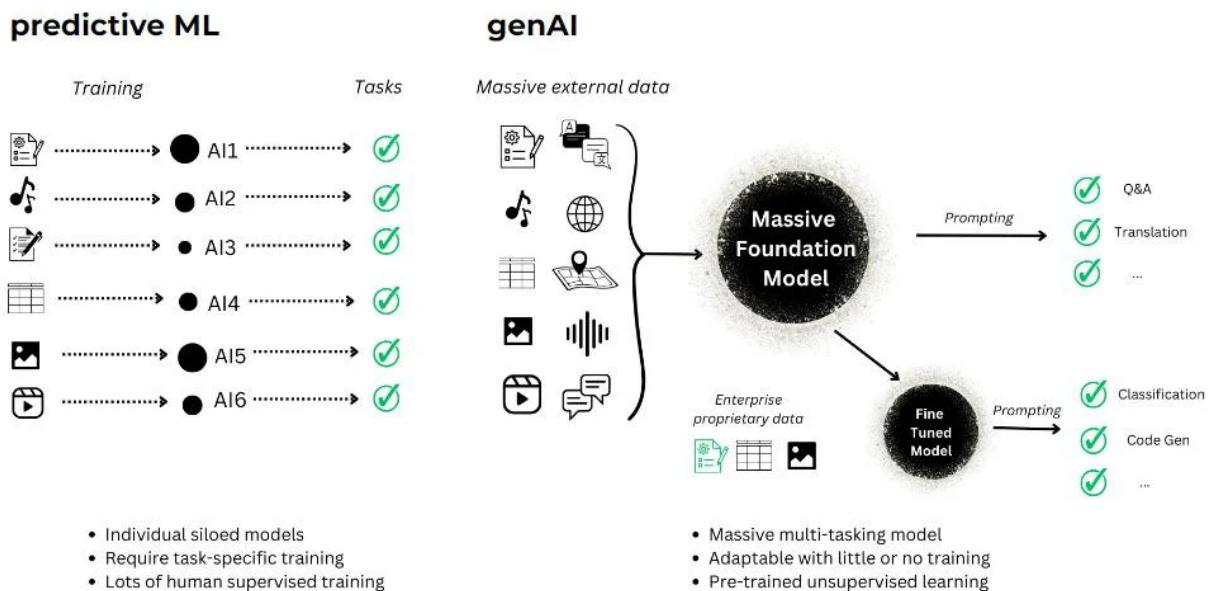
- **Application Examples:** Predictive analytics in business forecasting, spam filters in email systems, and recommendation systems in streaming services.
- **Key Characteristics:**
 - **Supervised Learning:** Often relies on labeled data sets to train models. Labeling data is very expensive and time-consuming.
 - **Predictive Accuracy:** Emphasizes the accuracy of predictions based on known data.
 - **Analytical Approach:** Aims to understand data and draw conclusions.

Generative AI

In contrast, Generative AI doesn't just analyze data; it creates new data that didn't exist before. It's about innovation and creation, generating new content that is similar to but distinct from the training data.

- **Application Examples:** Creating new images or artwork, generating realistic human-like text, or composing music.
- **Key Characteristics:**
 - **Creative Output:** Produces new content, extending beyond analysis.
 - **Model Types:** Uses models like GANs (Generative Adversarial Networks) and VAEs (Variational Autoencoders) for content generation.
 - **Innovation Focused:** Pushes the boundaries of what machines can create.

traditional ML vs generative AI



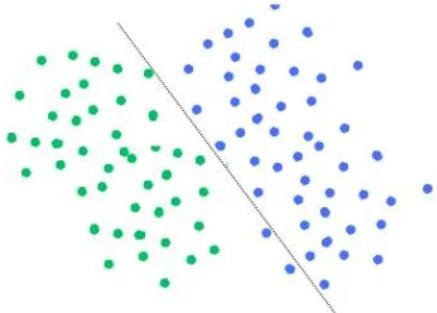
The Role of Discriminative vs Generative in AI

The distinction between discriminative and generative models is vital. Discriminative models excel in classification and prediction tasks, making them suitable for analytical applications. In contrast, generative models are unparalleled in their ability to create and innovate, making them ideal for tasks requiring new content generation.

deep learning model types

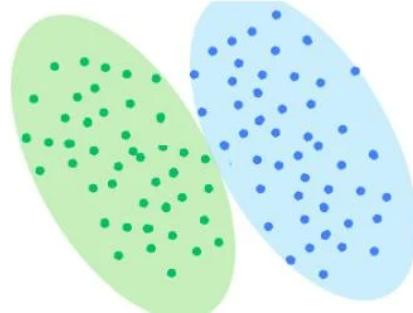
Discriminative

- Used to classify or predict
- Typically trained on a labeled dataset
- Learns the relationship between the features of the labeled data points



Generative

- Generates new data that is similar to data it was trained on
- Understands distribution of data and how likely a given example is
- Predict next word in a sequence



Understanding whether your business needs to analyze and classify existing data or generate new, unseen content will guide you in choosing the right AI approach.

Tomorrow, we'll explore another exciting aspect of AI. Stay tuned!

Best,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768, Seattle, WA 98104-2205

Subject: Day 2: Types of Generative AI Models
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Friday, January 5, 2024 2:50:46 AM GMT+05:30
Date Received: Friday, January 5, 2024 2:50:47 AM GMT+05:30

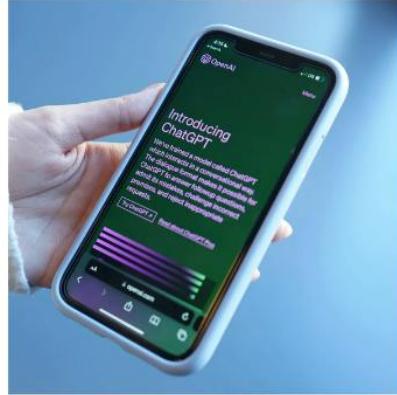
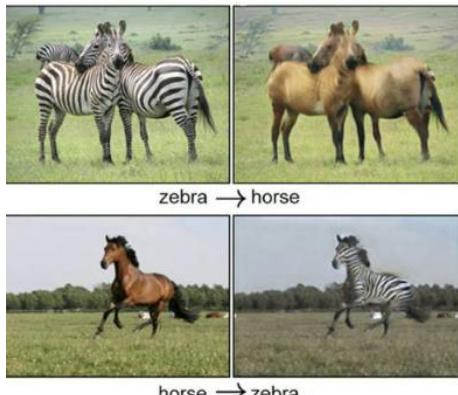
Day 2 - Types of Generative AI Models

Welcome to Day 2 of our Generative AI journey! Today, we're diving into the different types of Generative AI models, each with its unique capabilities and applications.

To keep it simple, I summarized this in four types:

- 1. Generative Adversarial Networks (GANs):** These models are game-changers in image generation, creating everything from art to realistic photos.
- 2. Variational Autoencoders (VAEs):** VAEs are great for tasks that involve compressing and generating high-quality images, offering applications in style transfer and more.
- 3. Transformer Models:** Known for their prowess in text, transformer models like GPT are revolutionizing text generation, translation, and automated writing.
- 4. Restricted Boltzmann Machines (RBMs):** RBMs excel in understanding complex data patterns, aiding in tasks like feature learning and topic modeling.

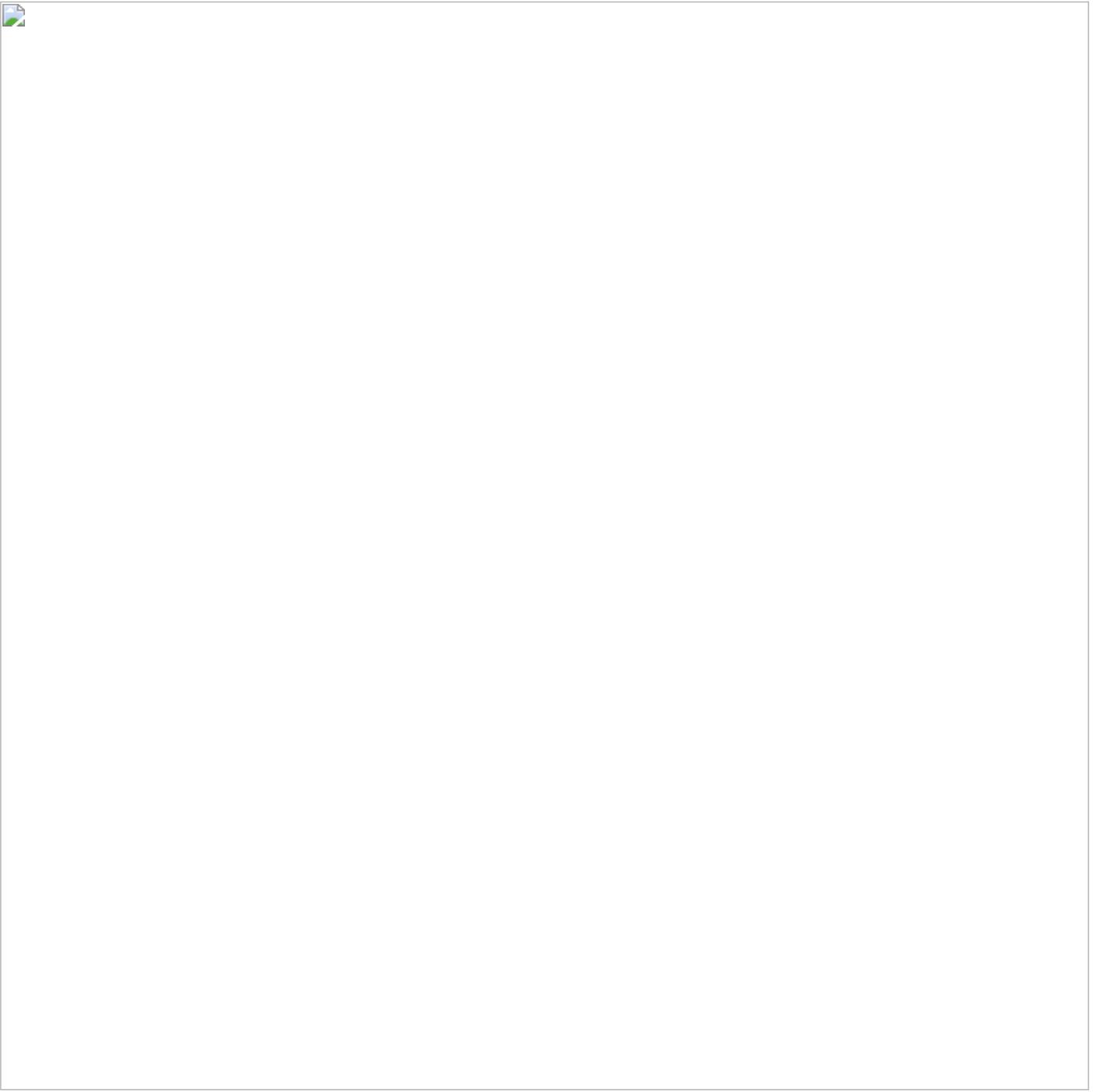
types of generative models



Generative Adversarial Networks

Autoregressive Models

Diffusion Models



But to simplify it even further, let's talk about Generative Language Models and Generative Image Models:

types of generative ai models

Generative language models

Generative languages models learn about patterns in language through training data.

Then, given some text, they predict **what comes next**.

Generative image models

Generative image models produce new images using techniques like diffusion.

Then, given a prompt or related imagery, they **transform random noise into images** or generate images from prompts.

- Generative language models learn about patterns in language through training data. Then, given some text, they predict what comes next.
- Generative image models produce new images using techniques like diffusion. Then, given a prompt or related imagery, they transform random noise into images or generate images from prompts.

There's a research paper that changed everything, called 'Attention is All You Need'.

The paper that changed everything

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin* †
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best

This paper introduced:

1/ Transformers: The paper presented the transformer model, moving away from traditional Deep Learning methodologies that were quite limiting such as Recurrent Neural Network (RNNs) and Convolutional Neural Network (CNNs) in NLP.

2/ Self-Attention Mechanism: Transformers use self-attention to efficiently process different parts of input data.

3/ Improved Parallelization: Transformers enable more efficient training through better parallelization compared to RNNs.

4/ Enhanced NLP Performance: They significantly outperform previous models in tasks like machine translation and text summarization.

5/ Basis for Advanced Models: The transformer architecture underpins major NLP models like BERT and GPT, enhancing language processing capabilities.

Here's the in case you want to go deeper, highly recommended to read: <https://lnkd.in/gaJCDKQH>

While I've focused on text and image generation, it's exciting to note that similar principles are being applied to **audio** and **video** generation, which is likely going to start exploding this 2024 and beyond. AI is now creating music, sound effects, and even generating or altering video content. The potential in these areas is vast and still unfolding.

and that's all for today. I hope you learned something insightful!

See you tomorrow!

Best,

Armand

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

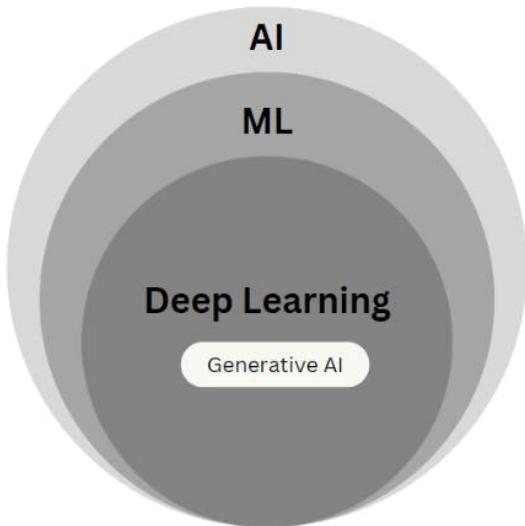
Subject: Day 1: Introduction to Generative AI
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Wednesday, January 3, 2024 2:47:59 PM GMT+05:30
Date Received: Wednesday, January 3, 2024 2:48:01 PM GMT+05:30

Welcome to Day 1 of our Generative AI journey!

Today, I'm uncovering what Generative AI is and why it's a game-changer in the business world. Imagine AI not just analyzing data but creating new, innovative content – that's Generative AI!

Ok let's start with the basics, but don't worry, we will get into more advanced concepts as we go.

definition



Definition of terms

Generative AI is a type of Artificial Intelligence that creates new content based on what it has learned from existing content.

When given a **prompt**, the model predicts what an expected response might be, creating new, original data like images, text, audio, video.

"Creativity" powered by examining large training datasets.

- **Artificial Intelligence (AI):** AI is the broad field of computer science focused on creating machines capable of performing tasks that typically require human intelligence.
- **Machine Learning (ML):** ML is a subset of AI involving algorithms and statistical models that enable computers to improve their performance on a task through experience.
- **Deep Learning:** Deep Learning is a subset of ML based on artificial neural networks, where algorithms learn from large amounts of data to identify patterns and make decisions.
- **Generative AI:** Generative AI refers to AI technologies that can generate new content, ideas, or data that are coherent and plausible, often resembling human-generated

outputs.

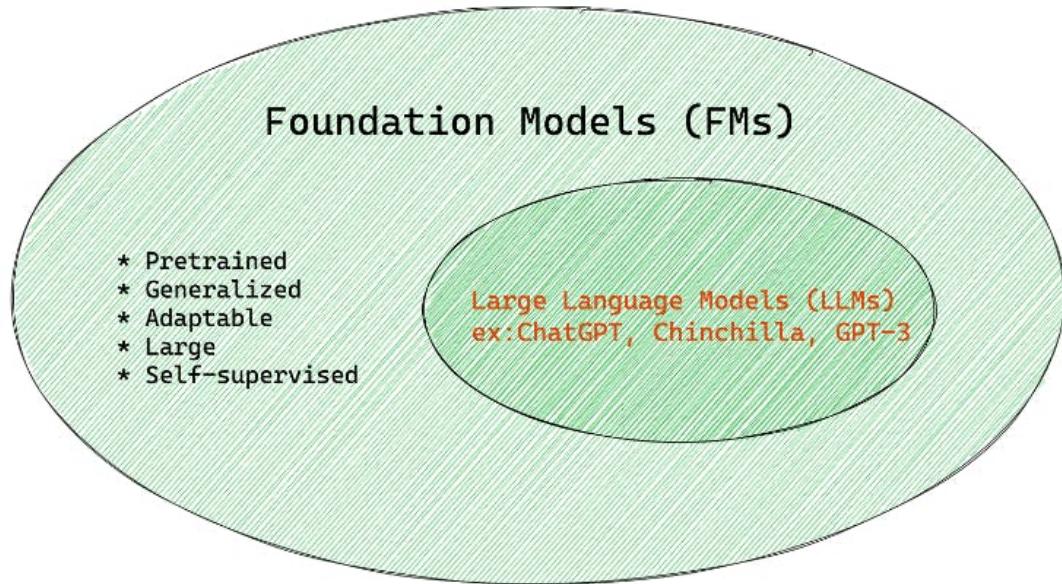
What powers Generative AI

Foundation models are large-scale artificial intelligence models that have been trained on vast amounts of data. These models are highly versatile and can be adapted to a wide range of tasks and applications.

Generative AI is one of the applications of foundation models. It involves using these models to create new content, such as text, images, or music. The foundation model serves as the underlying structure that understands and processes information, enabling the generative AI to produce new, coherent, and relevant outputs.

In simple terms, foundation models are like the core engine, and generative AI is one of the many things that this engine can power.

the models powering generative ai



FMs are models trained on broad data (using self-supervision at scale) that can be adapted to a wide range of downstream tasks.

What makes Foundation Models so powerful?

1. **Pretrained:** The model has already been trained on a vast dataset before being fine-tuned or applied to specific tasks.
2. **Generalized:** The model is capable of performing well across a wide range of tasks, not just the ones it was specifically trained for.

3. **Adaptable:** The model can be easily modified or fine-tuned to suit particular needs or tasks.
4. **Large:** The model is built with a substantial architecture and trained on extensive data, giving it a broad understanding and capability.
5. **Self-supervised:** The model primarily learns by analyzing and making sense of unlabeled data, without explicit guidance on what to learn.

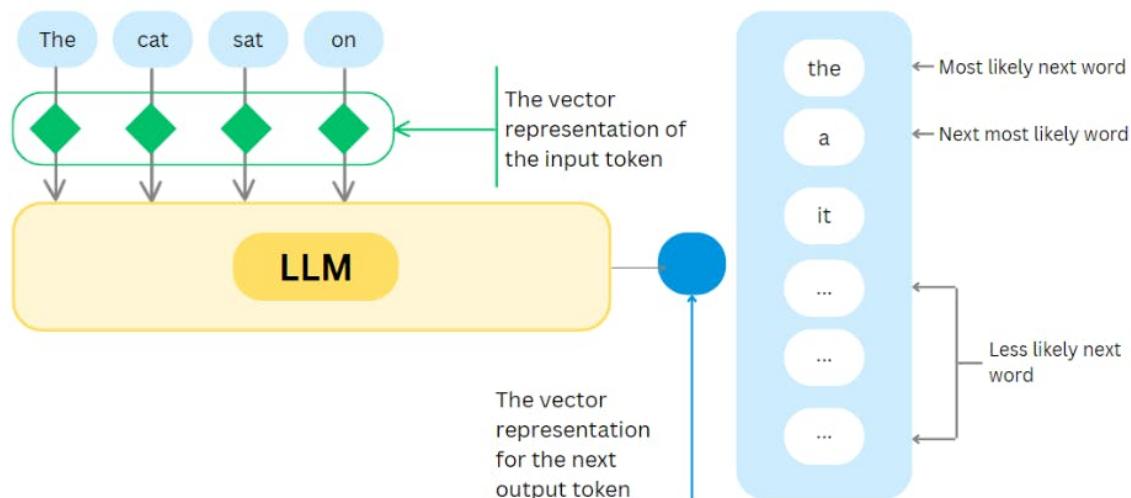
and what are Large Language Models?

Large Language Models (LLMs) are a type of foundation model specifically designed to understand and generate text. They're trained on huge amounts of text, which makes them good at a wide range of language tasks. LLMs are part of the broader category of foundation models, meaning they're versatile and can be adapted for different uses involving language.

LLMs like GPT take, as input, an entire sequence of words, and predicts which word is most likely to come next. They perform that prediction of the next word in a sequence by analyzing patterns in vast amounts of text data.

a next word predictor

LLMs like GPT take, as input, an entire sequence of words, and predicts which word is most likely to come next.



There's a big debate that LLMs do more than predict the next word; they compress a "world-model" within their complex networks and weights. This is an area of active debate within the AI community. You can join the discussion about this [here](#)

Two important concepts to understand in LLMs are:

- **Weights:** Numerical values within a machine learning model that are adjusted during training to influence the model's output in response to input data.

- **Parameters:** The broader set of configurable elements in a model, including weights, that determine its behavior and performance.
 - **Tokenization:** The process of converting text into smaller units (tokens), such as words or subwords, which are used as the input for LLMs to understand and generate language.
-

Ok's! That's it for today. I told you I would be short and right to the point.

Stay tuned for tomorrow's insights.

Best,

Armand 😊

PS. If you enjoy this free course, share it with your friends and colleagues! Here's the link:

<https://go.nocode.ai/30-day-email-course>

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205

Subject: Welcome to Generative AI in 15 days!
From: "armand@nocode.ai" <armand@nocode.ai>
To: deepakchawla35@gmail.com
Date Sent: Tuesday, January 2, 2024 6:38:30 PM GMT+05:30
Date Received: Tuesday, January 2, 2024 6:38:31 PM GMT+05:30

Welcome to the Generative AI email course

Welcome and thank you for enrolling in my **15-Day Email Course** on Generative AI! I'm thrilled to have you join this exciting journey where I'll unravel the potentials of Generative AI, specially tailored for business practitioners like you.

What to expect

Over the next 15 days, you'll receive daily emails, each designed to be a bite-sized, yet comprehensive exploration into various facets of Generative AI. From the basics of machine learning to ethical considerations and real-world business applications, this course is structured to enrich your understanding and spark your creativity in leveraging AI for your business needs.

- **Day 1: Introduction to Generative AI** - Overview of generative AI and its importance in business.
- **Day 2: Types of Generative AI Models** - Exploring different generative AI models like GANs, VAEs, and transformers.
- **Day 3: Traditional ML vs Generative AI** - Comparing traditional machine learning with generative AI methods.
- **Day 4: What are GPUs** - Understanding the role of GPUs in AI and machine learning tasks.
- **Day 5: What it Takes to Train a Foundation Model** - Insights into the resources and processes for training large foundation models.
- **Day 6: How to Customize Foundation Models** - Discussing techniques for customizing foundation models for specific uses.
- **Day 7: The Most Popular LLMs Available** - Overview of the most widely-used large language models and their features.
- **Day 8: Generative AI Applications and Use Cases** - Exploring practical applications of generative AI across business sectors.
- **Day 9: The Generative AI Stack** - Understanding the components and architecture of the generative AI tech stack.
- **Day 10: The Emergence of Small Language Models** - Discussing the rise and importance of small language models in AI.
- **Day 11: The AI Engineer Profession and Skills** - Exploring the role, responsibilities, and required skills of AI engineers.

- **Day 12: Ethical Considerations in AI** - Discussing the ethical challenges in AI development and deployment.
- **Day 13: Create Your AI for Business Roadmap** - How to develop a strategic AI integration roadmap for businesses.
- **Day 14: Future Trends in AI** - Exploring future developments and trends in AI.
- **Day 15: Continuing Your AI Journey** - Providing resources and advice for continued AI learning and exploration.

What's next

Stay tuned for your **first email tomorrow**, where I'll dive into the intriguing world of Generative AI. Get ready to embark on a journey that promises to enhance your understanding and reshape the way you think about AI in business.

Remember, each email is crafted to be a less-than-5-minute read, ensuring you can seamlessly integrate this learning experience into your busy schedule.

If you have any questions or need assistance at any point, feel free to reach out to us. We're here to support your learning journey!

[Unsubscribe](#) | [Update your profile](#) | 113 Cherry St #92768,, Seattle, WA 98104-2205