# AINA

### Magazine
AI and Analytics

Meet
# Josh Starmer
The founder of "Statquest"

## Urban Analytics
How AI can help in creating smart cities of the future

―――――

## Synthetic Data
Artificially generated data to train the next generation of AI models

―――――

## Privacy in ML
Techniques to preserve privacy while maintaining model performance

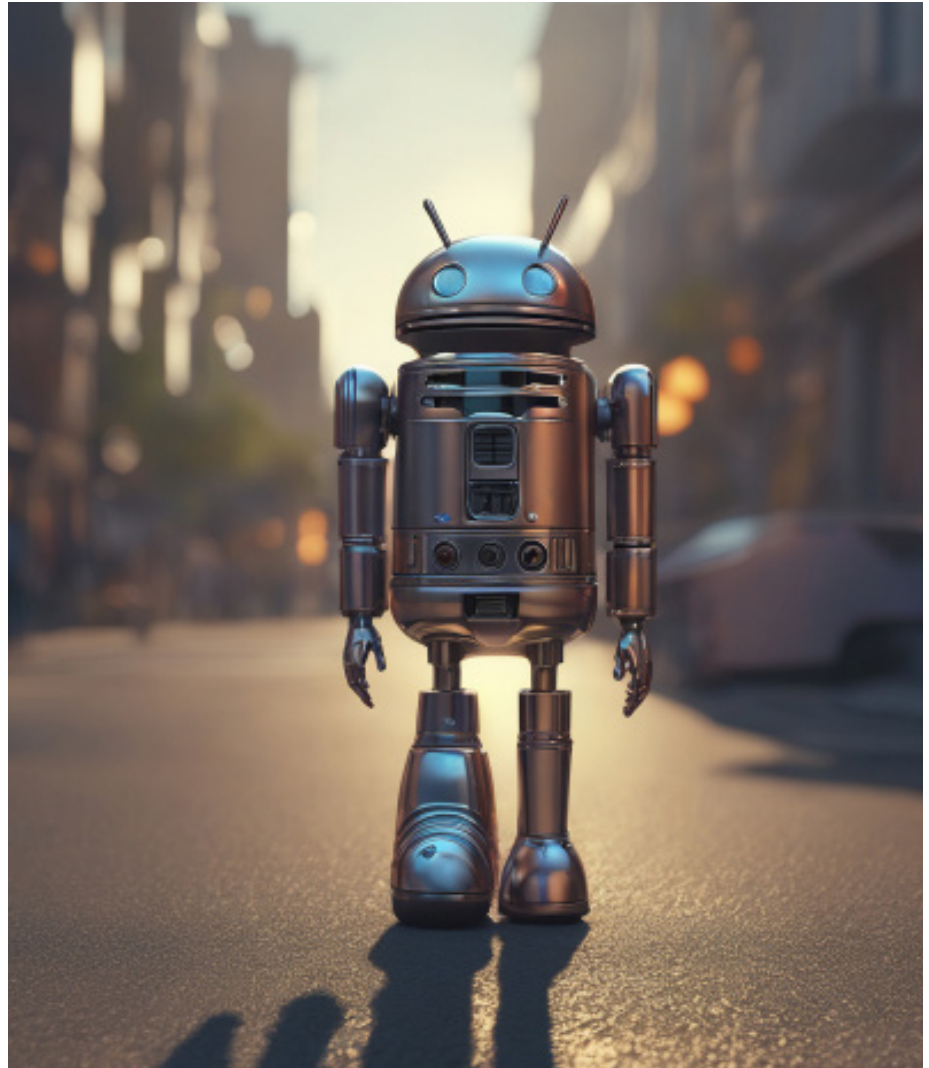# The Creative RenAIssance
## How LLMs are Shaping Generative AI

Writabrata Bhattacharya

# AI's Uncharted Voyage

Back in 1983, Claude Shannon, the renowned American mathematician and computer scientist, who is credited with inventing Information Theory and introducing Shannon's Entropy, expressed a fascinating vision for the future: Today, we stand on the verge of turning that vision into reality, witnessing a remarkable surge in the utilization of artificial intelligence over the past decade. From its origins as a purely academic pursuit, AI has rapidly evolved into a powerful force driving various industries and profoundly impacting the daily lives of millions of individuals.

In the world of technology and artificial intelligence, we've recently achieved remarkable milestones, developing AI systems that can learn from vast amounts of data, ranging from thousands to millions of examples. These cutting-edge advancements have revolutionized our understanding of the world, unlocking creative solutions to intricate problems. Thanks to these large-scale models, we now witness AI-powered systems effortlessly grasping



*"I visualize a time when we will be to robots what dogs are to humans, and I'm rooting for the machines"*

*- Claude Shannon*

human speech and written language. Think of the natural-language processing and understanding programs that have become a part of our daily lives, like digital assistants and speech-to-text applications. Beyond that, these extensive datasets, containing everything from the masterpieces of renowned artists to the entire collection of existing chemistry textbooks, have opened new frontiers in generative models. These models are now capable of producing awe-inspiring artworks inspired by specific styles or proposing novel chemical compounds based on the history of scientific research. The possibilities seem limitless as AI continues to redefine what's achievable in the realm of entertainment and technology.

In the realm of AI, while numerous systems are making a real impact on practical issues, their development and implementation demand substantial time and resources. To address specific tasks effectively, having a comprehensive and well-labelled dataset is crucial. Otherwise, human effort is required to painstakingly locate and label suitable images, text, or graphs. The AI model must then learn from this dataset before it can be tailored to your unique needs, be it language comprehension or discovering new drug molecules. Additionally, training a large-scale natural-language processing model can have a significant environmental impact equivalent to running five cars throughout their lifespan.

# Rise of Foundation Models

The next phase of AI is poised to revolutionize the dominance of task-specific models that have prevailed thus far. The future lies in the realm of foundation models, which are trained on extensive sets of unlabelled data, enabling them to perform diverse tasks with minimal fine-tuning. The concept of foundation models was popularized by the Stanford Institute for Human-Centered Artificial Intelligence through a comprehensive 214-page paper published in the summer of 2021.
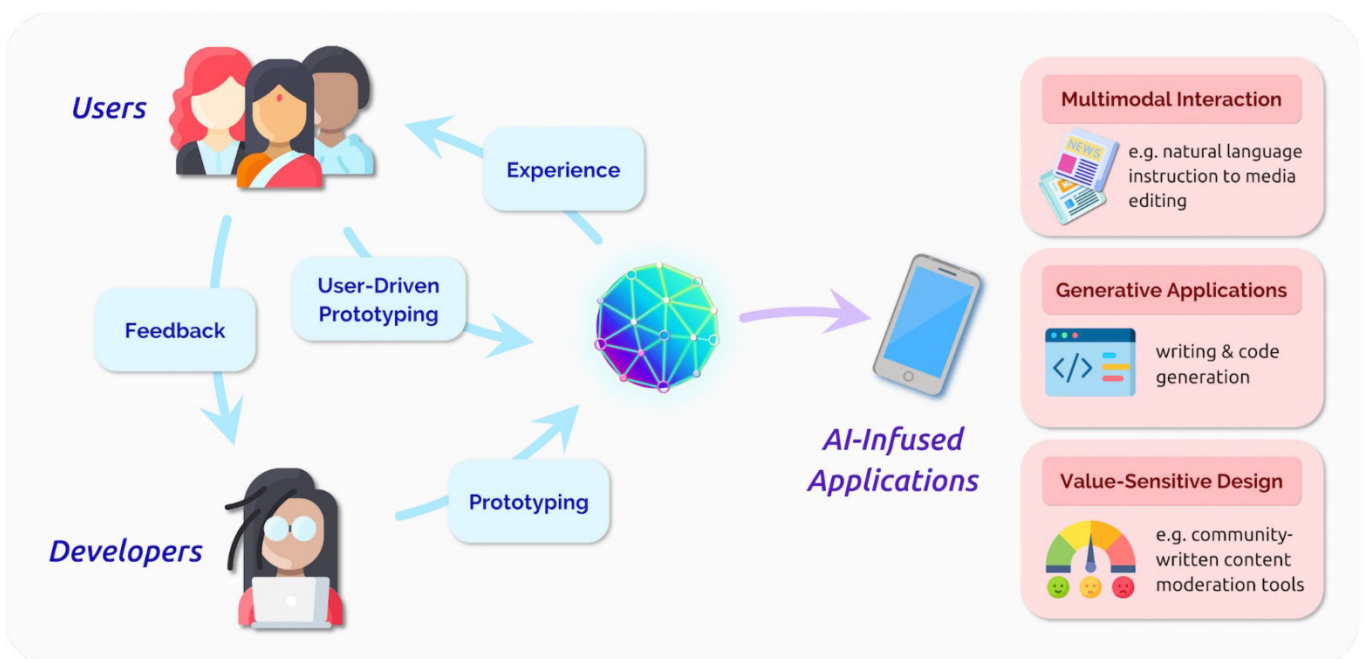
Early glimpses of the potential of foundation models have been witnessed in domains like imagery and language, with pioneering examples such as GPT-3, BERT, and DALL-E 2 showcasing their impressive capabilities. These systems can generate complete essays or intricate images based on brief prompts, even though they were not explicitly trained for those exact tasks or image generation in that precise manner.

Through self-supervised learning and transfer learning techniques, the model can apply the knowledge gained from one situation to another, much like how learning to drive one type of car allows for an easy transition to driving other types of vehicles with minimal effort. This advancement significantly boosts the efficiency and versatility of AI systems, enabling them to tackle new challenges without extensive retraining or reliance on specialized models for each task. As a result, the horizon of possibilities for AI applications expands, pushing the boundaries of what can be achieved with this transformative technology.

"The question of whether a computer can think is no more interesting than the question of whether a submarine can swim"

- Edsger W. Dijkstra

# Generative AI:
# Unleashing Creativity

The emergence of foundation models has paved the way for a transformative branch of artificial intelligence known as Generative AI. This incredible technology enables AI systems to create new and original content in different forms, like text, images, audio, and even synthetic data. Just like you draw or write stories, generative AI can do that too, but it's incredibly fast and can produce a multitude of amazing creations in a short span of time. It's like having a super-smart creative assistant that helps us explore fresh ideas and bring forth fun and exciting possibilities in the world!
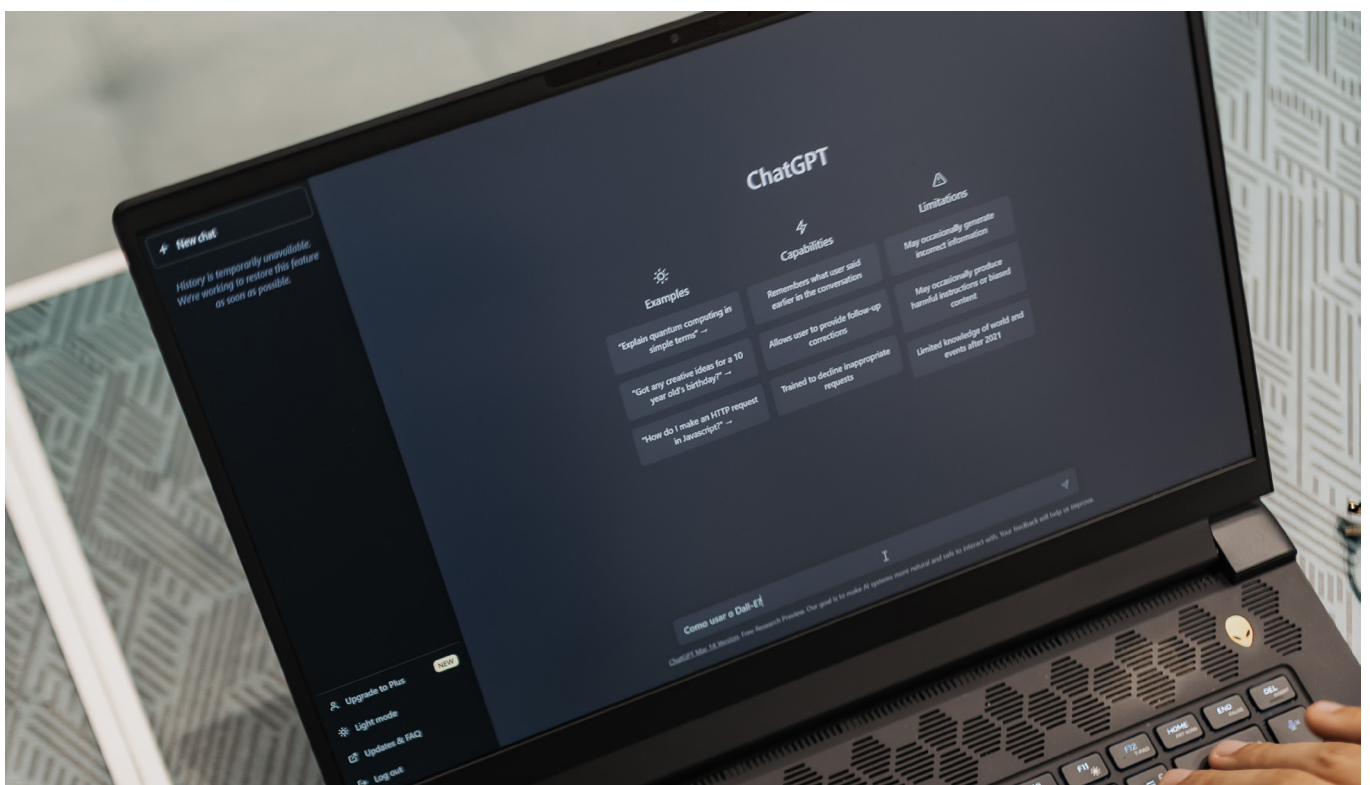
In contrast to other AI approaches that primarily focus on recognition, prediction, or classification tasks, Generative AI centres its attention on creating new data by leveraging patterns and structures learned from training data. In other words, it abstracts the underlying pattern related to input and uses that to generate similar content.

To illustrate this concept, consider a generative AI model trained on a dataset of human faces. This model has the capacity to produce lifelike and original faces that bear resemblance to the examples it was exposed to during training. Similarly, a gen

erative AI model trained on a collection of poems can generate new poems characterized by a similar style or theme, showcasing its creative potential.

Generative AI holds significant implications across a multitude of fields, spanning art, entertainment, design, and scientific research. In the realm of video games, it can be employed to generate realistic graphics that enhance players' immersion and visual experience. Additionally, generative AI can generate synthetic data that aids in training other AI models, facilitating the development and advancement of AI technologies. Moreover, it can serve as a valuable tool in creative endeavours, assisting in writing or music composition tasks by generating novel ideas or compositions. In the realm of scientific research, generative AI can contribute to drug discovery efforts by generating new molecular structures with potential pharmaceutical applications.

In recent times, two notable breakthroughs in the field of AI have captured significant attention: Large Language Models (LLMs) such as ChatGPT, Bard, and PaLM, as well as Audio Generation Models (AudioLMs) like MusicLM.



Generative AI enables AI systems to create original content in different forms in a short span of time
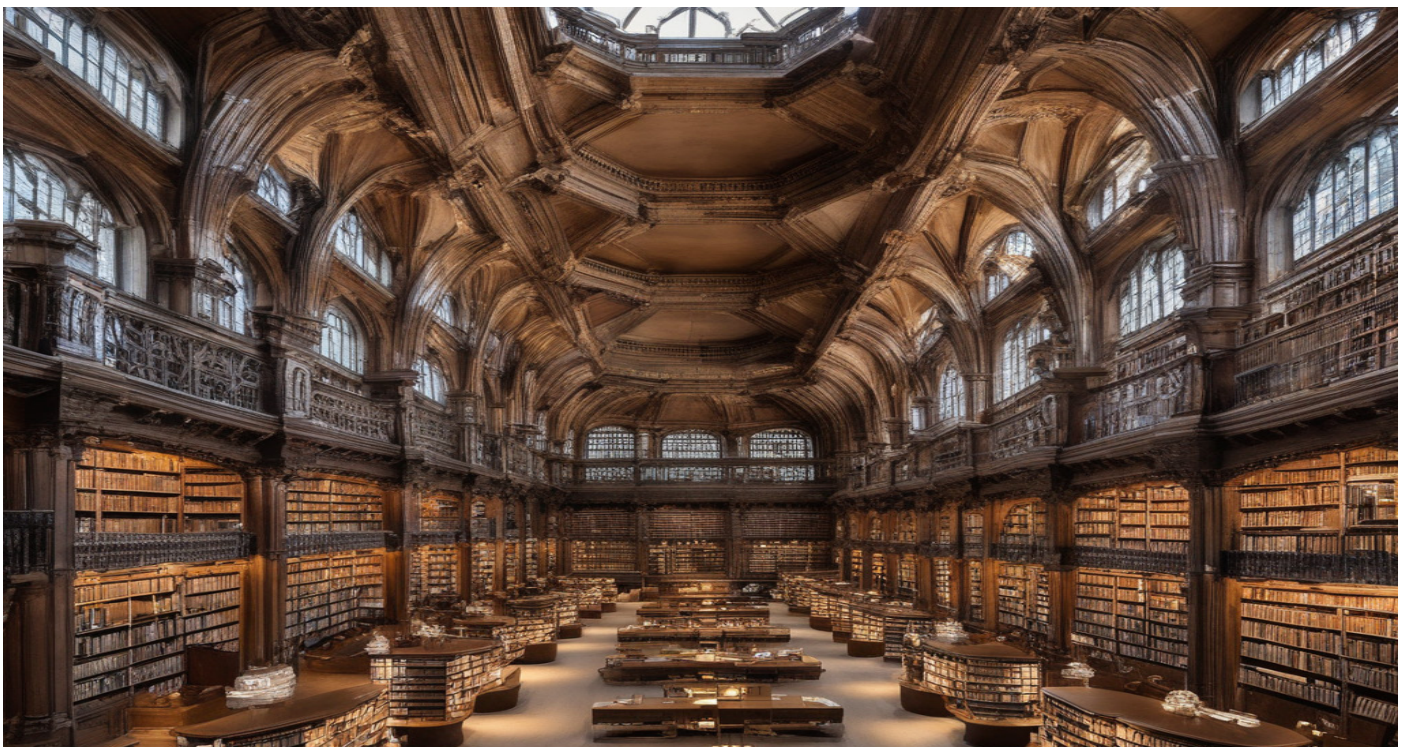
# Large Language Model

A large language model is a trained deep-learning model that understands and generates text in a human-like fashion. It is like a very smart robot that knows a lot about words and can help us with reading and writing. It has learned from many books and stories, so it understands how sentences are made and can answer questions or even create new stories for us. It's like having a clever friend who can talk with us and help us learn new things about language and words.

These models can be understood by breaking down their key concepts into three components: "Large", "General purpose", and "Pretrained & Fine-tuned".

datasets and a tremendous number of parameters. "Pretrained and fine-tuned" refers to the process of initially pretraining an LLM on a vast general-purpose dataset and subsequently fine-tuning it for specific objectives using a much smaller dataset.

There are several benefits to utilizing large language models. Firstly, a single model can be applied to various tasks. These LLMs, trained on petabytes of data and boasting billions of parameters, possess the intelligence to tackle tasks like language translation, sentence completion, and text classification. Secondly, LLMs require minimal domain-specific



LLMs trained on petabytes of text data possess the intelligence to tackle tasks like language translation

"Large" refers to the vast size of the training dataset, sometimes reaching the scale of petabytes, and it also pertains to the parameter count, which represents the memories and knowledge the machine acquires during model training. Parameters determine the model's proficiency in solving specific problems, such as text prediction.

These models are called "General purpose" as they are versatile enough to address common problems. This concept arises from the universality of human language, which transcends specific tasks, combined with resource availability constraints. Only a select few organizations possess the capability to train such large language models with massive

training data when tailoring them to address specific problems. Even with limited training data, LLMs can perform reasonably well and exhibit the ability to recognize patterns that were not explicitly taught during training. Lastly, the performance of LLMs continues to improve as more data and parameters are added.

Over the past few months, LLMs have left a strong impression in various fields. From crafting poetry to helping with vacation plans, we're witnessing impressive progress in AI's abilities, creating value for businesses and individuals alike. It's even quite possible that the very article you are reading right now was generated by AI itself.

# Pathway Language Model

One prominent example of a Large Language Model is PaLM, developed by Google. PaLM, short for Pathways Language Model, is a dense decoder-only transformer model that leverages the innovative pathway system, enabling Google to train a single model across multiple TPUs. This model exhibits proficiency in various domains, including mathematics, coding, advanced reasoning, and multilingual tasks such as translation.

During its training, PaLM was exposed to a diverse range of data, encompassing scientific and mathematical information, 100 spoken languages, and over 20 programming languages. It serves as the underlying technology for Google's workspace products, Bard, and the PaLM API. It is possible that we have already been utilizing PaLM without even realizing it. Notably, PaLM is not limited to translating spoken languages but can also handle programming languages.

The versatility of PaLM is exemplified by its ability to facilitate collaboration between individuals with different language backgrounds. For instance, an English-speaking individual can use PaLM to collaborate with a colleague in a codebase where all the documentation is written in Korean. Moreover, PaLM excels in generating and comprehending nuanced language, including idioms and riddles. This capability is crucial as it requires understanding not only the figurative meaning of words but also the literal intent behind them.

As LLMs continue to evolve and be integrated into various applications, we can expect further advancements in language processing and understanding, benefiting individuals and organizations across diverse linguistic and professional contexts.



PaLM 2 is a language model developed by Google that aims to bring AI smarts to some of Google's most popular apps, such as Gmail, Google Docs and Bard.

# Music Language Model (MusicLM)

*"Musik kann die Welt verändern"*
*- Ludwig Beethoven*

As Ludwig van Beethoven famously stated, which translates to "Music can change the world." When it comes to music generation, there is indeed a remarkable potential for transformation. One notable model in this realm is MusicLM, which has the ability to generate music from text captions without utilizing any diffusion techniques. What sets MusicLM apart is its ability to consistently generate music at 24 kHz, maintaining a high level of fidelity over several minutes.

MusicLM is a cutting-edge audio model that operates using two types of tokens: semantic and acoustic tokens. The semantic tokens capture the essence of melody and rhythm, while the acoustic tokens, provide valuable information about the recording conditions of a track. These tokens come together in a remarkable fusion, resulting in the creation of high-quality audio tracks. One of the most impressive features of MusicLM lies in its melody conditioning capabilities. This means it can generate music that precisely matches a given text prompt while incorporating the desired melody seamlessly. What's even more astonishing is that MusicLM has been put to the test with descriptions of paintings, showcasing its remarkable ability to translate the mood and atmosphere of visual art into mesmerizing musical compositions.

Another intriguing feature of MusicLM is its generational diversity. The same text prompt can yield a range of music compositions or variations within the same sample. This flexibility opens up countless creative possibilities for musicians and composers. AI has the potential to revolutionize music production in the future, envisioning a scenario where a smart keyboard can generate numerous arrangements with different instruments, sounds, and styles based on a simple melody played by the musician. Instead of engaging in the laborious process of composing and producing music from scratch, artists can guide AI systems to create captivating tracks until the desired result is achieved.

The integration of continuation and voice cloning techniques further expands the horizons of music production. Imagine the ability to bring legendary artists like Michael Jackson and Kurt Cobain back to life, enabling them to sing a duet on your track using their cloned voices. Additionally, leveraging ChatGPT, you can employ AI assistance to craft the lyrics for your compositions. To complete the immersive musical experience, Vision AI can be utilized to capture the mood and emotions displayed on the faces of the audience, allowing real-time generation of music that perfectly suits their reactions and preferences.

As advancements continue, musicians, producers, and listeners can expect a groundbreaking era of creativity, collaboration, and personalized musical experiences. The power of music to shape the world is undeniable, and AI's pivotal role in its evolution promises a thrilling and harmonious journey ahead.



MusicLM can create high quality music from text captions

# Conclusion

In conclusion, the emergence of Generative AI marks an exhilarating new era filled with endless creativity and boundless potential. From the transformative power of Large Language Models like PaLM to the enchanting melodies crafted by MusicLM, AI systems are reshaping industries and revolutionizing artistic expression. As inclusivity, guiding us towards an AI-driven entertainment landscape that reflects the best of humanity's ingenuity and empathy. The future of generative AI is a symphony of human ingenuity and machine intelligence, harmonizing to create a world where imagination knows no bounds. By embracing AI's potential with a mindful approach, we can shape a

*"The real problem is not whether machines think but whether men do"*
*- B.F. Skinner*

we continue to push the boundaries of what AI can achieve, the possibilities are thrilling but uncertain. As we explore the vast potential of generative AI, it becomes imperative to confront the complexities of responsible deployment. Our creative endeavours must uphold the core values of privacy, fairness, and future where creativity flourishes while safeguarding against unintended consequences. Let us embark on this journey with optimism, curiosity, and a deep commitment to using AI ethically and responsibly, shaping a world where art, music, and beyond thrive with the best of human and machine collaboration.



Music Boom in Greece. 100% Available Elvis Presley

## UDAI SHANKAR
## Director- Data Science, Providence

*Udai Shankar is a data science leader and researcher having 18+ years of experience in solving real-world business problems in various domains including operational intelligence, NLP and security. Udai has previously held senior roles at CA Technologies. He has strong knowledge of algorithms with experience in understanding and shaping use cases.*

**To start can you share about your personal journey in data science and machine learning?**

I started learning data science sometime in 2011. Prior to that, I was a C++ developer with CA technologies, and we've done some interesting work on compilers and a lot of good stuff on C++ core programming stuff. Then in 2011 I took a course in pattern recognition from IIIT Hyderabad, and that's how I got started. Around 2013 I shifted totally into data science with CA technologies and joined Providence in 2019.

**Your interests also lie in topological data analysis and quantum machine learning. These are very niche topics involving mathematics and computer science. What inspired you to pursue research in those specialized areas?**

I love mathematics, that's what got me interested in machine learning in the first place. I studied as a part-time student at IIIT Hyderabad. During this time, I found quantum computing very interesting. I studied theoretical computer science and then quantum computing at IIIT itself. I am an independent researcher in Quantum Algorithms, more specifically, I study the homology groups that arise in topological data analysis from the lens of quantum hidden subgroup problem.

**That's fascinating. So I myself come from a mathematical background from CMI. We had algebraic topology in our curriculum in both undergrad and master's. In the theory of homology, wherein we can count the number of holes in some topological spaces. I want to know how this whole detection formalism of algebraic topology can be connected with these robust computational methods to extract qualitative features in data.**

The area that I research is an intersection of computational complexity, quantum hidden subgroup algorithms and computational topology. Since you are into math, you might have studied, algebraic topology, which cannot be learned without the algebraic fundamentals – basically abstract algebra (groups, rings and modules). There is a result in quantum computation that says that if there is an abelian group, it is easy for a quantum computer to compute the generators of a hidden subgroup in that abelian group. This result is important because all the quantum algorithms where an exponential speedup was observed can be reduced to instances of abelian hidden subgroup problem (AHSP). For example the Shor's algorithm, discrete log, period finding etc., are all instances of AHSP. There are also some results when the groups are not abelian [like Normal subgroup, dihedral subgroup etc.]. Now, (free) abelian groups arise naturally when we work with homology. Topological data analysis (TDA) links data and topology (persistent homology etc.). So I study the groups that arise in TDA from a quantum lens. There is one famous result on quantum algorithms for TDA by Seth Lloyd, but this is not from the algebraic perspective (i.e., not directly using AHSP).

So the way I entered into this is just out of curiosity and exploration. Math excites me, so machine learning excites me, and this also excites me. Topological data analysis is very powerful. Not the best out-of-the-box algorithm, but probably the internals of it can be used in your day-to-day data science as well. That's my perspective.

"Math excites me, so Machine Learning excites me"

**Okay, that's actually nice to hear because, from what we have seen in theoretical sciences, it's difficult to connect with reality if you go into those niche topics. So are there any applications of topological data analysis that have been going on in the industry?**

There is an entire company that spun out of topological data analysis called Ayasdi. They started with topological data analysis. There is one problem with topological data analysis. Our edge node graphs are one-dimensional topological objects. In machine learning, we have probabilistic graphs where each node is represented with a probability distribution, so that's how we study machine learning, isn't it? Everything in machine learning is about an implicit graph & (implicit) probability distribution on the nodes, and once, we have the probability distribution, we are trying to predict something using the probability distribution. So basically, under every scenario, we have a probabilistic graph. Now, think about a scenario where blood pressure is, let's say, one node and diabetes is another node. We might have an edge between them, that is a 1D graph. Suppose, we ask, how do blood pressure and diabetes function together? When a person has both diabetes and blood pressure, what is the curve that determines their, let's say, disease progression? In other words, if I am studying BP, Diabetes and the future state of CKD together, I have to take not just the edge between diabetes and blood pressure, I have to take all three together, and that's a simplex. So when you take three nodes together, then the entire triangle is a

simplex. That's how higher-dimensional topological objects enter into the picture. An interesting area to explore is persistent homology.

**AI algorithms have already matched the performance of human experts on several prediction tasks. But humans still have some valuable domain knowledge, which hasn't been incorporated into the learning process. In this context, how do healthcare professionals collaborate with AI systems to combine human expertise in decision-making?**

I'll answer this from first a theoretical perspective and then come to the practical aspects of it. So, from

an ML algorithm is to provide a causal graph [or a higher dimensional topological object]. Your algorithms work on the causal graph with the support of the data.

**As an extension of the previous question, how do you see the role of Bayesian models in healthcare AI and how do they differ from traditional machine learning models? Could you give some examples of healthcare applications where Bayesian models have really shown advantages?**

The main power of Bayesian modelling is creating custom distributions and inferring them using MCMC or variational inference, right? For example, let's say, mod-

## "Science is about discovering the causal links , I propose a causal link and verify whether it is causing it or not"

the theoretical perspective, there is science, and then we have all our machine-learning algorithms that can predict, cluster etc. Suppose you're trying to do a controlled experiment, where you want to see whether a particular drug really affects a particular disease. The domain knowledge must be passed as an input to the algorithm because domain knowledge is not available in data, so it has to be passed as an input to the algorithm. Data can also be used to validate different domain perspectives. Think about it this way. How do I present my domain knowledge to the algorithm? It must be a mathematical object. So, my take on this is, that the best way to present domain knowledge to

el a coin. I know that coin can be modelled as a Bernoulli distribution. Now the question really is, suppose I want to model a different type of coin – maybe a coin with a memory. That means it remembers its previous toss. Suppose I claim I created this coin and give it to you. Now model this coin in terms of a probability distribution. It's way easier and clearer in a Bayesian setting than in a non-Bayesian setting. Also, the causal graphs and all that I was talking about, can be easily correlated or linked to a Bayesian setting than a non–Bayesian setting. If I know that F causes A, I can directly put a distribution on F and A and then link them together, and the problem becomes extremely simple to solve. So, Bayesian mod-

elling is very, very fundamental. It provides you with tools to model real-world problems. Suppose I were to give the same information to a deep learning algorithm without any Bayesian modelling [Of course, Bayesian deep learning & Graph NNs exist, but for the moment, imagine we are not using any of it] I'm actually not telling the algorithm anything about the causal links (the Science or the Domain) or about the probabilistic graph that is available from the domain. The causal structure is better represented naturally through Bayesian learning. So, for non-Bayesian approaches, you have to do something unnatural about it.

**Now that we have a lot of talk about data privacy and there is a Data Protection bill in India. So in the context of healthcare, how do you handle sensitive or confidential patient data while communicating visualizations to your stakeholders?**

Yeah, so there are two levels to it. One is where we work within the internal structure of providence, where we have access to data. Before I use data for a specific purpose, there's a governance board known as the IRB Board whose permission we take and then work on our research project. From an external perspective, when I want to give out data to, let's say, a partner. We de-identify the data (HIPPA Compliant data). We make sure that all the PHI and PII data is totally masked (or obfuscated), and we change dates in such a way that the deltas remain the same [so that ML is still possible on the data]. We also mask out the zip codes, which have less than, let's say, 20,000 people. The ages are converted into buckets so that age groups with few people cannot be used to identify the person. These are the kinds of things we do for privacy. All the data that goes to any algorithm usually goes through this de-identification process. We have a service that does de-identification on the data.

**As a follow-up question, is there any regulatory compliance which you have to follow in case you have to take data of US patients and work in India?**

It's just the same, we go through the IRB process, as I said, which is the internal regulatory board. Providence India is just a part of Providence overall, so there is no separation or anything out there.

**As you earlier said, like all of us are hooked to ChatGPT and large language models these days. So**

in light of the availability of advanced language models such as ChatGPT, modern chatbots are being used to enhance customer experience in the case of healthcare. Considering there have been instances where chatGPT is giving out wrong answers and giving out inaccurate information, how is the healthcare industry making sure that it can leverage these models to its advantage?

We obviously have to finetune (or retrain) these models for our purposes. As it is, the hugging face LLM models are not very useful to us because they are huge and because of the high computational requirements. Also, the open-source LLM models are not trained for the specific use cases we are interested in. Let's say you are a doctor and you have a specific question in your mind and you come to the clinical data (EHR data) to gather some output as to what the data is saying about your Hypothesis (the specific causal question in your mind). So you've got a causal link that you want to validate through the data. That's one type of question. The other type of question is, let's say, for example, how many people are coming to the emergency department, i.e., the footfall at the emergency department, and how

"Chat-GPT is not a plug-in, but the inherent technology is very useful still. You know, there is a dialogue element to it, a human feedback element to it."

to schedule nurses and things like that.

Let's talk about yet another kind of problem where I want to predict the revenue of the company. So I want to look at how the revenue is getting generated and where I can improve and get more revenue and all that. That's another kind of scenario. None of these scenarios are directly amenable to working with the LLMs. The data is usually a collection of disparate pieces of structured and unstructured information that need to be combined and linked together in a systematic way in order to run any deep learning algorithm on it.

**You were talking about, extending this to images and genomics data, for LLMs. Could you talk about what exactly you do there when you're using images or genomics data for LLMs?**

We have not started working on the images and genomics yet. Even with the EHR data, there are two parallel worlds if you think about it. The unstructured data is one world of information, for example, patient notes - the doctor has written some notes, and the images have their diagnostic reports. Then you have this structured data, which includes things like the BP measurement or, let's say, their heart rate, and pulse rate. The structured data also includes the medication, diagnosis etc. The first problem really is how do we bring these two worlds together. My claim is through discretization, vectorization and graph NNs. For example, we can build vectorial representation from textual data on structured and unstructured data. We can bring out a

connection.

between words that are close by. Then I can draw an edge between them. So, I made it into a graph, thus shifting from the textual to the graphical picture. If we want to do a proper analysis of the data, our LLMs have to work over a causal structure derived from data and domain knowledge.

**From your answer, it was clear that ChatGPT can't be the plug-in, plug-out solution for everything. Is it?**

Exactly. It is not a plug-in, but the inherent technology is still very useful. The backend, which is, Transformers + RLHF is extremely powerful. But this has to be made to work over the data specific to our domains and with an understanding of the causal structures inherent in my problem.