**40+**
**YEARS**
**SINCE 1982**

# DATAQUEST

*CyberMedia*

THE BUSINESS OF INFOTECH

# GENAI:
# WALLFLOWERS
# ARE NOW
# WALL-FACERS

It's a 4 body problem with GenAI.
It's so tough to get that syzygy right-
where data availability, data privacy,
authenticity and enterprise-readiness
align to the T and not cause chaos.
But how tough, exactly?

# GenAI: Wallflowers are now Wall-facers

It's a 4 body problem with GenAI. It's so tough to get that syzygy right- where data availability, data privacy, authenticity and enterprise-readiness align to the T and not cause chaos. But how tough, exactly?
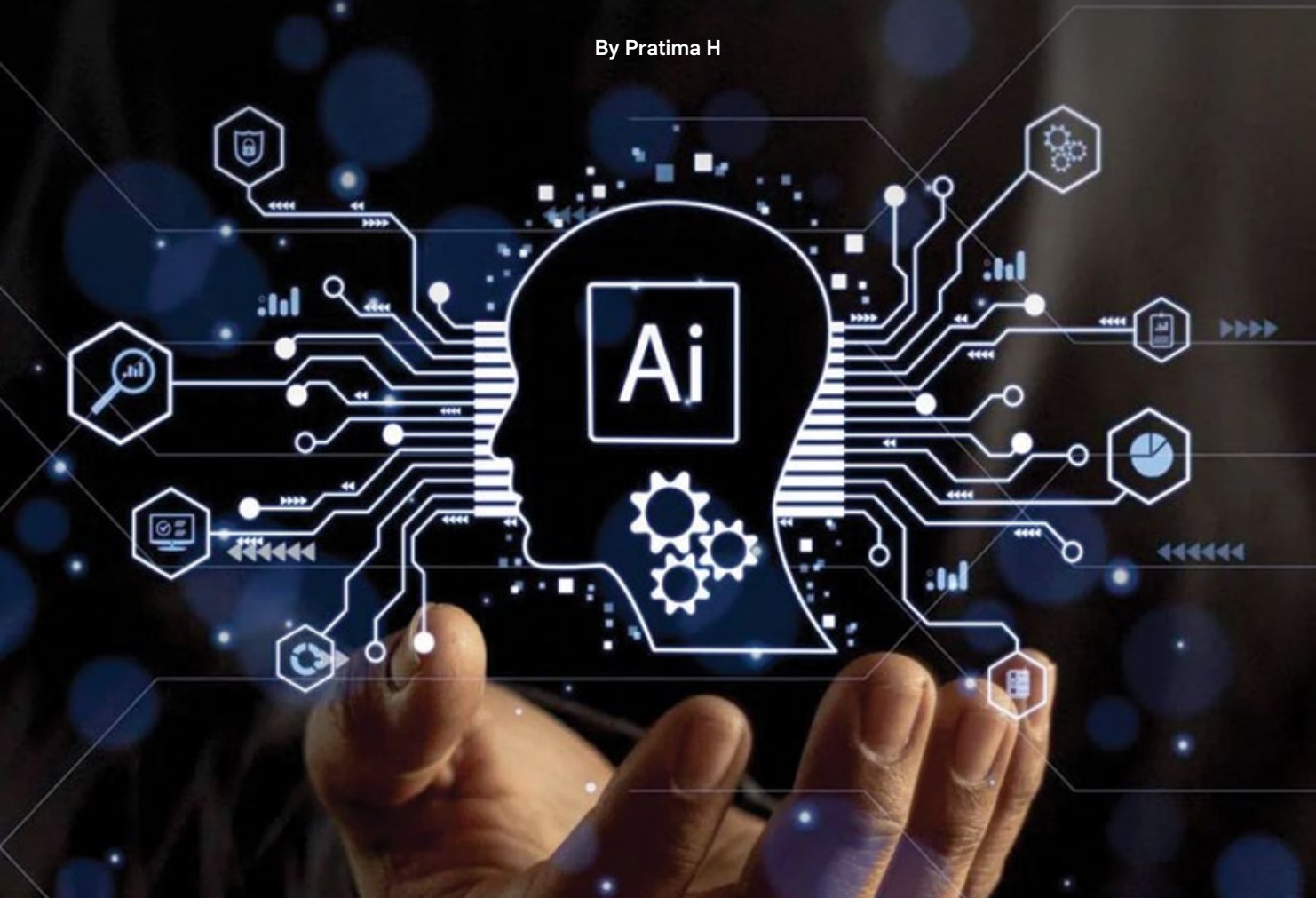
**By Pratima H**

*"Overly literal translations, far from being faithful, actually distort meaning by obscuring sense."*
— Ken Liu, The Three-Body Problem

"Liars." The super-intelligent and much-advanced-than-humans species from the other world is thrown off by just this one word. In 'The Three-Body Problem' series, the dark plans of aliens (for lack of a better word) to invade this planet go off the orbit just because they hear a story. They listen to a human as he casually narrates a fairy tale about the Wolf and Red Riding Hood. And they get stuck at the word, the idea, beyond rescue. They just don't understand it. Humans can say one thing

and mean something else. And the other person will actually 'get' it.

That's exactly the power and potential of GenAI. It's the most beautiful, complex and rarefied accomplishment of AI. Because it's about human language. It's about words, thoughts, nuances and context. That takes a whole new weight when we are thinking of a country like India- where vocal search, speech-based apps and mobile phones (in the hands of people unable to read and write) are as rife as rich literature, multiple languages and text waiting to be digitized. GenAI has to 'get' what humans say here. And if it can do that, there's no stopping it from ruling this world.

This is where the classical challenge of data availability in GenAI assumes an entirely new complication. It has to dance well with other factors like data privacy, data authenticity, context and enterprise usability. Not that easy a dance though.

Let's translate them one by one and find out why.

## LOCAL PLUS AMPLE - A TRICKY DATA UNIVERSE

It's a massive white space – waiting for the right trajectories to be formed. From about $67 billion in 2024, the global GenAI market is slated to grow to an explosive $967 billion by 2032 (as per Fortune Insights). Other crystal balls don't have such fiery graphs. But still potent enough. Grand View Research pegs this space at around $109 billion by 2030 and Allied Market Research puts it at about $191 billion by 2032. If we go by what McKinsey estimated in 2023, generative AI could add the equivalent of $2.6 trillion to $4.4 trillion annually across the 63 use cases it analyzed. But there can be many a slip between this cup and the graph- in India.

In the global GenAI landscape, the US and China are leading the way as they are creating their own GenAI models built with informational and contextual safeguards. However, India seems to lag in the fray, as the strength of GenAI, in the context

of local languages, relies on expertise and models trained specifically for Indian data points- points out Sanjoy Paul, Program Director, Tech and Data, Hero Vired. "Thus, unlike Western solutions, which may not be as effective due to cultural and linguistic differences, there's a need for AI models to be trained on Indian datasets."

Large Language Models (LLMs) are sophisticated prediction machines that are good at finding patterns amongst vast and diverse set of data that they are trained on, Lakshminarasimman Raghavan, Group VP Technology, Publicis Sapient gives a quick primer first. "They use this to generate new data/content based on their training. Data is not just limited to text, but these models can also generate images, videos etc. As far as data exchange with LLMs for GenAI use cases goes, two types of data are involved, one, the data that is passed to LLMs as prompts and the data that the LLMs are trained on. There is a third consideration stemming out of this - which is how the data supplied as prompts are used by the LLMs to train their base models."

In any application landscape, data availability is critical, a principle especially resonant in the specialization of GenAI, as Sridhar Mantha, CEO of GenAI Business Unit, Happiest Minds Technologies adds next. The challenges – as he rightly underlines- are two-fold: content and context.

The first one is- How is the data being made available (Content) – Static information in the form of documents ranging from static documents in various formats to multimedia, alongside transactional data accessed through APIs, databases, etc. Each source poses its unique intricacies and demands in terms of integration and interpretation. The second is - What is the data (Context) – How do we understand the context of the data is what presents the crux of the challenge, particularly evident in transactional information. Understanding the context requires a nuanced



**"**

In India, the primary challenge is data regulation. For instance, standards such as DPDP have come into effect only last year and the need of the hour is strict enforcement of these privacy laws.

**- Lakshminarasimman Raghavan,** Group VP Technology, Publicis Sapient

> **"**
> By addressing both content and context with equal emphasis, organizations can harness the full potential of their data assets.
> **- Sridhar Mantha,** Happiest Minds Technologies

grasp of the intricacies surrounding each data point—contextual cues, temporal relevance, and relational dependencies. For GenAI applications, this contextual understanding is fundamental for accurate analysis, inference, and decision-making," Mantha explains.

Indic LLMs will also face data scarcity as a challenge and hence have a reduced amount of data to train their models on, Raghavan contends. "Given this situation, the possibility of ambiguous responses will be a problem for such models as more such ambiguous responses are likely to reduce their authenticity."

Ramke Ramakrishnan, VP Analyst, KI Leader at Gartner brings in another critical facet here. "The availability of data influences the level of processing applied to the data sets, which could contribute to hallucinations, the accuracy of the results and the effectiveness of security controls. These factors, amongst others, can also impact the overall costs, deployment time and expertise required to deliver a robust GenAI application."

Maybe one way to counter this availability problem is to create synthetic data.

With the growing adoption of synthetic data across industries and vendors, we see an immense value and increase in the adoption of synthetic data to improve ML performance, fairness and accuracy with more training data sets and fine-tuning models and privacy preservation on sensitive data to deliver better outcomes. Also, it allows organizations to take on new use-cases for which they need more real data." Shares Ramakrishnan.

As to synthetic data, Ramasubramaniam Srinivasan, Senior VP - Practice Leader - AI/ML, Data and Cloud of Collabera Digital reminds that sourcing of authentic datasets for fine-tuning is going to be crucial to avoid copyright-related concerns. "Also there is a real paucity of local data sets and synthetic data is estimated to be 40-50 percent of the data used for development of AI models. Several start-ups have started offering tools, platforms and services for synthetic data generation."

Gaurav Kheterpal, founder and CEO of Vanshiv Technologies opines, "While India is a goldmine for data from various sources, I believe synthetic data is also important for LLMs due to factors such as regional dialects, cultural nuances, the strong need for domain-specific content and adoption to societal context."

Synthetic data is expected to become 100 percent of the data used for AI model development in a few years, Srinivasan ventures to augur. "Synthetic data may reduce the problems regarding personal information sharing and promote data minimization. However, there are also deep-fakes issues which have been recently witnessed in India due to synthetic data generation."

That's exactly where we move to the next challenge of GenAI in India.

**PRIVACY AND SAFETY - ECLIPSES WORTH WATCHING**
Ok, let's say we got the data we need – and in the numbers we want. What next? The Scrabble gets harder from here.



> **"**
> With AI algorithms processing amounts of data, verifying the accuracy and reliability of outcomes becomes essential. In industries such as finance and online retail where AI plays a role in decision-making, ensuring the accuracy of data inputs and transparency in algorithm operations is essential.
> **- Murale Narayanan,** Dell Technologies

**"**

Establishing accountability for lapses is crucial, prompting discussions around the necessity of human validation or the use of additional AI systems for validation.

**- Srinivasulu Nasam,** Bosch Global Software Technologies

We also have to ensure that user privacy and safety are not compromised. And issues like bias, prejudice, offense and misinformation are not allowed to breed in the rush for more and more data. Data privacy is a very delicate, but crucial, challenge in GenAI.

According to the latest Gartner Business Outcomes of Technology by Use Case Survey, senior leadership's main concerns with GenAI are 'Privacy concerns' and 'risk of misuse'. Despite the potential for AI to deliver value, there is a deluge of distrust that needs to be resolved to make it purposeful, Ramakrishnan shines the torch a little closer on this issue. "With AI-related regulations increasing, adopting open-source foundation models is preferred over closed-source due to deployment flexibility, customization options and enabling better control over security and privacy."

Data privacy is a concern, particularly due to the abundance of data generated and handled by AI systems, seconds Murale Narayanan, CTO, Managed Services India (Technology Transformation), Dell Technologies. "Practical examples vividly demonstrate the complexities involved. For example, in healthcare, where AI systems analyze data, it is essential to implement robust measures, for data privacy to protect patient confidentiality. Using encryption protocols, access controls and anonymization techniques can help mitigate risks and maintain privacy standards."

Dinesh Kumar Poobalan, CEO & CTO, Greatify, a venture into digital learning and education, also avers that data privacy and authenticity are

paramount in GenAI models, as they determine the trustworthiness and reliability of the technology. "Data privacy is also important in education because in order for individuals to be willing to engage online, they have to trust that their personal data will be handled with care."

Currently, models are trained using data available on the internet, which may include copyrighted material in various formats, articles, and content produced by individuals, as well as certain personal data protected by relevant data protection guidelines- cites Srinivasulu Nasam, Senior Technical Director – GenAI for Systems and Software Engineering at Bosch Global Software Technologies. "Regulations, usage scenarios, and potentially compensation models need to be defined for the ethical use of such data in model training."

Arun Moral, Managing Director at Primus Partners illustrates how the Indian Government continues to promote GenAI through 'Make AI in India and Make AI work for India' announced in 2023 budget with fiscal grants to COE setup for GenAI. For example, the passport verification process has embraced GenAI techniques for swift validation of an individual's identity. The process leverages the Machine Readable Zone (MRZ) code approach to accelerate and enhance overall operational efficiency. The accuracy of the GenAI tool to recommend or reject the application is directly proportional to the quality and correctness of the base data which is a huge underlining threat.

Moral also reminds of the concerns arising from unsanctioned GenAI tools in enterprises. His



**"**

Synthetic data is a new class widely used in various aspects, including data anonymization, AI and ML model development, data sharing and monetization. Synthetic data is at the peak of inflated expectations from the most recent Gartner Hype Cycle for Generative AI.

**- Ramke Ramakrishnan,** Gartner

> **"**
> "Large Indian organizations have matured significantly in ML implementation as it has been ongoing for a few years. This is expected to continue."
> **- Ramasubramaniam Srinivasan,** Collabera Digital

prescription to the challenge of Shadow IT here is to embrace the magic pill envisaged in 'Fifth Industry Revolution's - Human-centred-Architecture'. "This would ensure industry outcome complemented with accelerated shoes of GenAI and controlled by the human brain."

Rahul Bhattacharya, Technology Consulting and AI leader at EY GDS notes that it is likely that existing LLMs (both proprietary and open source) have been trained on a certain amount of material related to India or produced within. "As the range, quantity, quality, and sources of this data remain unknown, it's hard to determine the authenticity of this data, or whether it contained private data that was used without permission. Any error in training due to inauthentic or biased data could lead to models that are ineffective or, worse, perpetuate stereotypes and inaccuracies."

Mantha points out how Enterprise GenAI applications primarily rely on data sourced from within the enterprise itself. "These applications predominantly serve to interpret and visualize this internal information, leveraging GenAI technologies to derive insights and facilitate decision-making processes. To maintain the integrity of this data boundary and uphold privacy standards, numerous technical interventions are deployed. These interventions encompass a spectrum of measures, ranging from robust encryption protocols and access controls to anonymization techniques and data masking strategies."

Sanjoy Paul, Program Director, Tech and Data, Hero Vired affirms that to meet legal obligations and defend user rights, preserving the privacy of personal data used to train AI models and produce content is indispensable. The Personal Data Protection Bill, which strongly emphasizes permission, data localization, and severe penalties for non-compliance, presents hurdles for Indian markets. At the same time, preserving authenticity in generated material is essential to stop fraud, disinformation, and manipulation of deep fakes.

There is also the question of the protection of Enterprise IP. Given the computational requirements and the nature of LLMs, it may not be practical to run separate instances for each enterprise, leading to shared usage among multiple entities, Nasam cautions. "Architectural and design considerations should prioritize maintaining data security and privacy within each enterprise, with emerging concepts like RAG and enterprise-specific adaptive layers aiming to address these challenges."

**ENTERPRISE-READY GENAI - NOT IN THE ORBIT YET**

Just making models will not suffice unless there are takers and confidence in the scalability and relevance of such solutions. GenAI in India will have to find an enterprise-fit beyond experimental work.

Ganesh Gopalan, Co-founder, and CEO, Gnani. ai offers some examples. "We collaborate with numerous large banks to automate their customer support functions. Additionally, Gnani.ai assists over 50+ lending customers, such as Buy Now Pay Later companies, NBFCs, and microfinance institutions, in automating customer conversations throughout

> **"**
> While the adoption at the enterprise level is increasing, there's a need for further development in scaling the solutions and addressing local language contexts to ensure comprehensive integration across enterprises.
> **- Sanjoy Paul,** Hero Vired

> **"**
> Data privacy and authenticity are not just legal and ethical requirements but also critical business imperatives.
> **- Rahul Bhattacharya,** EY GDS

the lending life cycle." Moreover, they partner with automotive companies to facilitate early demand assessment and gather customer feedback through automated conversations with showroom visitors. "Furthermore, in the healthcare sector, we collaborate with global enterprise healthcare companies to automate patient engagement processes across various channels."

Let's ask someone who packs the knowledge of several pilots and implementations across government processes and banks in India. Ankush Sabharwal, CEO of CoRover, the conversational artificial intelligence (AI) startup behind BharatGPT, India's homegrown GenAI initiative raises an important point- the advantage of less time to go live and true representation of Indian context (how Indians ask and communicate) with enterprise-level solutions made in/for India. "They also take less time in training, for questions, for inferences and are apt for CRM and ERP responses. We have been doing virtual assistants for a long time now with Generative AI based on local LLM, a lot has improved – be it in the area of query resolution or go-live time windows of virtual assistants." He cites many pilots that are underway across a diverse spectrum – from a lot of banks to Chennai AI police to education solutions and e-government projects of some states like J&K.

"We are building Chatbots/ ML Solutions for use-cases here." Srinivasan from Collabera Digital, lets on. "We have a significant client base in India. Large Indian organizations have matured significantly in ML - and this is expected to continue. With Generative AI solutions starting from ChatGPT last year, there has been a significant interest in AI models in Small and Medium Organizations as well. With 'Pay as you go' models in Generative AI service providers, Indian organizations are currently at the POC stage and we see a lot of interest."

**ANSWERS - THEY ARE COMING**
To fully leverage AI in India, it's crucial to prioritize building LLMs for local languages using data generated by companies and government bodies, sums up Paul.

Nasam adds how we're already witnessing platforms like 'Bhashini' bringing this technology closer to a vast user base by overcoming language barriers. "However, key concerns such as data privacy, copyrights, intellectual property (IP) protection, and quality assurance must be addressed for generative AI to become mainstream."

The GenAI ecosystem is evolving to include small language models (SLMs), offering more streamlined operations than LLMs, Ramakrishnan offers an alternative. "SLM requires less data and fewer parameters for fine-tuning and training, making them promising to be more efficient and effective for specific domains or targeted applications. This shift improves data privacy control and enhances the user experience by reducing errors and increasing authenticity."

As Bhattacharya captures. "It is essential for AI solutions to not only be context-aware but also to have robust mechanisms that protect user data and reaffirm the reliability of the information provided."

Looks like there is a lot to do.

Not exactly light years away, but GenAI is still far from realizing its ultimate potential in a region like India. We have to cross many speed bumps around privacy, context, and scale to get to the future that awaits us. It's a thin line. But it can change the entire axis of this universe.

Spoiler Alert (If you haven't read or watched 'The 3 Body Problem'). It's as thin as that tiny difference between these two lines. "You are a bug." And "You are bugs."

Be it GenAI or the big universe out there - One small alphabet can change so much. 🅓🅠

*pratimah@cybermedia.co.in*