# The Nexus of Bitcoin and AI

By the Team at Spirit of Satoshi
December 2023

This report started as a summary of findings from the product discovery interviews we conducted in October, to help guide our product development process. It was supposed to be for internal use only. Along the way, however, we found many interviewees were eager to learn more. As a result, we felt it would be valuable to turn the internal report into a more substantial "industry report" for public consumption.

In the following pages, you'll find not only the data associated with our findings from the product discovery calls, but also insight into the process of building the Bitcoin-centric large language model, along with statistics on this nascent industry. The report also identifies some of the key challenges faced by bitcoin companies and entrepreneurs today, and highlights some of the potential applications, where a bitcoin-centric AI tool could add value or solve problems.

We've learned an incredible amount to date, and hope this report clears up a number of very common misconceptions about "training" models, tuning them, augmenting them, tooling, data curation, creation, cleaning, vectorizing, storing, quality assurance and so much more.

We hope you walk away from reading this with a deeper understanding of not only AI and Language Models, but with some clarity on what is actually going on here, what is real, what is useful, what is hype and some ideas about how AI can contribute to the growth of the bitcoin ecosystem.

Founded in 2023, Spirit of Satoshi is a Bitcoin-centric Language Model Project. The models are currently being trained on a carefully curated corpus of data that includes Bitcoin literature (books, essays, guides and podcasts), Austrian Economic and Libertarian literature, along with a small sprinkling of other related resources. People from across the global Bitcoin community are also contributing to the model's development by answering bitcoin-related questions, and helping verify the accuracy and relevance of units of data.

# FOREWORD

**F**ounded in 2023, Spirit of Satoshi is a Bitcoin-centric Language Model Project. The models are currently being trained on a carefully curated corpus of data that includes Bitcoin literature (books, essays, guides and podcasts), Austrian Economic and Libertarian literature, along with a small sprinkling of other related resources. People from across the global Bitcoin community are also contributing to the model's development by answering bitcoin-related questions, and helping verify the accuracy and relevance of units of data.

# EXECUTIVE SUMMARY

**T**he potential applications of a bitcoin-centric AI tool are wide-reaching. But the word potential here is key. Much of what is sold online about AI is hyperbole, and creates the illusion that it can do more than it really can. It reminds me of some inverse of the following image:

Something more like: "AI capabilities are less incredible and useful than they may appear."

This doesn't mean Language Models and other ML or AI tools won't or can't add significant value to the Bitcoin ecosystem (or any other industry for that matter). Much like the products we use today that leverage "AI", whether Uber and Google or your phone, there are obviously ways in which automation can be applied to scale up operations, speed up certain processes and make products and/or services better.

I use the word "automation" specifically here, because when it comes down to it, that's really what we're talking about. The big shift with LLMs is that we're now able to somewhat automate tasks - or elements of tasks - that require the use of language, or semantic reasoning.

It's still too early to say how much this will change the world, and whether it will have the size of impact that some say it will - but I am pretty confident that once the hype dies down, we will, over the coming decade, find clear applications and uses for such a tool.

In the meantime, join us as we analyze what is and is not useful by identifying the key challenges, and opportunities where Bitcoin overlaps with AI.

First some statistics. 280 participants have contributed over 40,000 responses in the fine-tuning process of the Spirit of Satoshi model, while over 33,000 bitcoin resources have been added to the Nakamoto Repository. We have by no means used all of this data for the training, but what we have used has been drawn from this pool, cleaned, formatted and used.

We interviewed a blend of Bitcoin companies, content creators and investors, totalling almost 50 and identified 5 common challenges across the board:

→ Marketing and customer acquisition,
→ User onboarding and effective education
→ Customer support
→ Employee onboarding, upskilling and technical development,
→ Hiring and scaling

Interviewees understood the value of having a truly differentiated, trained and fine-tuned language model that isn't captured by the mainstream or embedded with mainstream biases. Together, we identified product opportunities that cover a suite of different applications, including but not limited to:

→ Customer support & success agent
→ Bitcoin intelligence agent
→ Bitcoin content generation assistant
→ Bitcoin "tutor" or education assistant
→ Bitcoin "influencer"

Each of these products would be predicated upon the existence of an underlying "Bitcoin model" in order to operate effectively.

Beyond this data, we will examine the process of actually building a Bitcoin language model.

→ What does training mean, and how is it different from fine-tuning?
→ What is retrieval augmentation, and why is everyone using this as a means to develop "their own models"?
→ Why is such framing inaccurate, creating expectations that cannot today be met.
→ Why the primary costs for training a model are not, as most people assume, borne from GPU cycles, but from the data preparation stage.
→ Why the quantity of data has far less to do with the final product, than does the quality of the data.
→ How crowdsourcing can be used for the development of both general and domain specific models, whether open or closed source

How Bitcoin and Lightning can enable crowdsourcing at scale, and privately if needed.

In summary, this report is going to be full of valuable data, both quantitative and qualitative, new mental models and a whole lot of lessons. You will walk away understanding AI in far greater depth.

## ALEKSANDAR SVETSKI

# CONTENTS

# PART 1:
# DATA, FINDINGS AND COMPARISONS

**W**hile the Spirit of Satoshi model is built and tested, we sought insight into the current state of bitcoin businesses and the bitcoin landscape. Before the team develops and delivers a product or product suite leveraging a Bitcoin-centric LLM, we wanted to better understand the highest value opportunities for the future direction of Spirit of Satoshi.

We used a double diamond approach to product discovery, uncovering recurring challenges and converging on the main pain points. As we get further into the discovery cycle, we will continue to develop potential solutions that the Bitcoin-centric LLM can solve. Finally we will prioritize MVP product features to deliver built on top of the bitcoiner-trained Spirit of Satoshi.

The first section of the report will focus on the data we gathered, both in the interviews we conducted, and in the training process.
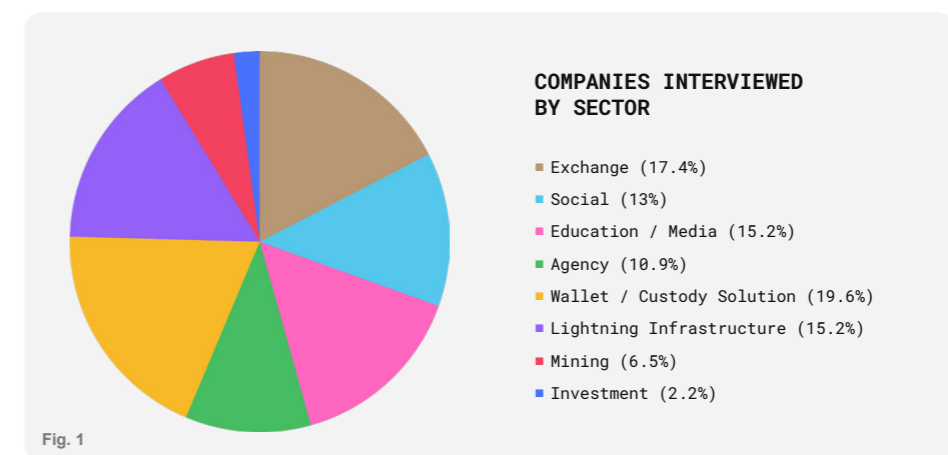
## METHODOLOGY

We connected with businesses, product/service-providers and content creators from across the Bitcoin and Lightning Network ecosystem. In doing so, we uncovered a series of common themes.

The product discovery exercise was based on a set of consistent interview questions, asked by our team to the interviewees in each of the sessions. Conversations were allowed to flow beyond the scope and explored relevant topics for each bitcoin business.

The team, consisting of Aleksandar Svetski, Alan Bakli and Jonathan Gordon, completed a total of 46 interview sessions across 45 companies, lasting between 30 minutes to one hour each. The 55 individual participants' fields ranged from bitcoin exchanges and wallets, finance and payment applications, social, review and educational platforms, bitcoin mining companies, media and Lightning Network infrastructure (see Fig. 1).

Interviews were conducted in September and October 2023 over Zoom and in-person, providing valuable insight into their business operations and future goals.



**COMPANIES INTERVIEWED BY SECTOR**

- ■ Exchange (17.4%)
- ■ Social (13%)
- ■ Education / Media (15.2%)
- ■ Agency (10.9%)
- ■ Wallet / Custody Solution (19.6%)
- ■ Lightning Infrastructure (15.2%)
- ■ Mining (6.5%)
- ■ Investment (2.2%)

Fig. 1

# 1. KEY FINDINGS

**B**efore we dive into AI solutions, the report will lay out the most common pain points companies are facing today in general, of which AI could potentially plug in as a helpful tool. The insights gleaned from the product discovery interviews offer a diverse range of perspectives from companies, entrepreneurs and content creators working on bitcoin, the Lightning Network, and the ecosystem at large.

Several common challenges emerged throughout the 46 calls we conducted. This section reports on the key findings, organized by pain points and opportunities. To protect anonymity and business integrity, no specific comments are attributed to any company.

The top challenges include customer acquisition, onboarding users along their bitcoin education journey, customer support and internal difficulties with bitcoin technical development, hiring and scaling.

Although more interview time was focused on the problem space, several opportunities were also suggested by participants and highlighted below.

## 1.1 PAIN POINTS AND NEEDS

Top challenges facing each company varied depending on the bitcoin business segment and size of the company. To begin the interview, each company provided the three biggest challenges facing their business today. We dive into each of the top challenge areas with additional responses from the question set (see Fig. 2).

### MARKETING AND SALES

→ 28 companies (64%) highlighted either customer acquisition, marketing, monetization, scaling, brand awareness and brand building among their top three challenges.
→ The majority of bitcoin businesses we interviewed are organically growing through word of mouth, in-person networking, referrals, podcasts and social media - primarily Twitter/X and LinkedIn. A handful of companies utilize influencer marketing and search engine optimization. Overall, paid advertising spend at the top of the funnel is low.
→ Bitcoin and Lightning companies are finding ways to earn consistent revenues, with many trying to reach "nocoiners" that are getting their first experience with bitcoin.

→ Content generation is time-consuming and needs consistent focus to build brand awareness with a distinct voice. Not all companies can afford this commitment. While some have played around with ChatGPT, it doesn't meet their expectations for bitcoin-related content. Most interviewees say that it's "good for ideation but not for anything I can use in public".
→ Monetization of bitcoin products and recurring revenues is a constant challenge, particularly in the bear market, while many incentivize usage and engagement by rewarding users with Sats.
→ 6 companies (14%) specifically mentioned copywriting support and media script writing as a top challenge. This falls under the marketing, particularly for content creators.

### ONBOARDING AND EDUCATION

→ 22 companies (50%) noted bitcoin education and onboarding as a top challenge. While each business is providing varied products and services, all bitcoin and Lightning companies are faced with introducing people to bitcoin.
→ Companies expressed the importance and challenge in onboarding nocoiners in an engaging way that gives them confidence.
→ Intuitive user experience (UX) and design is a critical component to bitcoin education and onboarding, which 3 companies (8%) specifically mentioned as a pain point.
→ A common need to communicate accurate educational content while not overwhelming the user was expressed. Companies need an improvement beyond simple links to FAQs, blog posts or googling to research on their own.
→ Companies are constantly dealing with general noise that creates confusion amongst users and recurring bitcoin FUD.
→ The quantity of solid bitcoin content has increased significantly since the last epoch, yet it can be challenging to feed the customer the right information. Focus on "meeting them where they are" or "aligning to their model of their world".
→ Users don't know how to learn more about bitcoin or ask the right questions. Supporting bitcoiners on their educational journey.

### CUSTOMER SUPPORT

→ 9 companies (20%) reported customer support as a top challenge.
→ Most companies provide limited, manual customer support mostly via email, or Telegram groups. Some companies use tools such as Zendesk or Intercom, which now have ways to build out automated responses to particular questions. A common occurrence is the snowballing of tickets when additional questions emerge before the ticket can be closed.
→ Amongst exchanges, the level of customer service can vary widely depending on the value of the customer, with service ranging from automated responses, to white glove.
→ There are common, recurring questions that users ask bitcoin onramps, offramps and exchanges. These include "where is my bitcoin transaction?" and other mempool-related questions, "how do I withdraw my bitcoin?", and questions related to the basic functions of a bitcoin or LN wallet.
  › These questions range between 40% and 80% of customer support efforts.
  › Companies are concerned with scam mitigation and helping users nervous about doing bitcoin or Lightning transactions.
  › A growing number of companies need to support customers in different languages, particularly those operating in Europe.

### TECHNICAL DOCUMENTATION AND CODE-WRITING

→ 8 companies (18%) highlighted coding, software development and technical documentation as a top challenge.
→ Given the decentralized and open source nature of Bitcoin, it is difficult for developers to all stay on top of new product developments in Bitcoin and the Lightning Network.
→ Even within companies, dependencies on internal documents are constantly changing, leading to a heavy burden for internal teams.
→ Current AI tools are not useful for Bitcoin and Bitcoin-related coding.

### EMPLOYEE HIRING AND TRAINING

→ 5 companies (11%) discussed the need for better tools to screen candidates, as well as training new employees on Bitcoin knowledge. Several noted the time wasted early on elements of the hiring process that could potentially be automated.
→ Although many are seeking to hire existing bitcoiners with the relevant skills, companies note that as the industry grows there will be a greater need to upskill for bitcoin.

### OTHER TOP CHALLENGES

→ Regulatory, legal & compliance —6 companies (14%)
  › Several are expanding into new jurisdictions. Expansion requires significant effort to understand local regulatory and legal frameworks for bitcoin. Especially relevant for exchanges, onramps and wallets.
→ Standard startup challenges (financing, time constraints) - 8 companies (18%)
  › Fundraising during the bitcoin bear market has been a challenge.
  › Bitcoin and Lightning companies are generally bootstrapped with individuals having limited time for redundant tasks, many of which can be supported with AI.

→ Internal communication & operations - 4 companies (10%)
  › Ability to retrieve information from internal documents and reference previous conversations had in messaging apps such as Slack.

It was impossible to note all of the challenges that came up in all of the conversations, but the above is a good cross section of what was common and recurring in the calls we conducted.

Following the initial part of the interview, we went on to discuss potential ideas and opportunities with many of the participants. These findings follow.

## 1.2 OPPORTUNITIES

Where problems lie, opportunities are to be found. Therefore, understanding the key challenges and identifying how to ameliorate them is where we decided to spend the balance of our time.

In this section we will outline some of the primary opportunities a product suite leveraging a Bitcoin LLM may be able to capitalize on. The participants provided creative responses for how the Spirit of Satoshi, either as a standalone model, or as part of a broader toolkit, could add value to their business. This was of particular interest, knowing that the balance of this decade will see a significant influx of people coming into Bitcoin and along with that will come demand for knowledge, education, support, tools, dev assistance and more.

It's important to note from the outset that the ways in which language models can be used and effectively integrated into a business's workflow are very new. We are all still figuring things out. Some ideas may seem obvious when you initially think about them, but there remains a large gap between theory and practice. Implementation remains a key challenge. The following list of ideas will follow the customer journey, from top of funnel, and on through the more internal development tools.

### BITCOIN CONTENT GENERATION AGENT

An agent focused on developing bitcoin-centric content for marketing and sales, including tweets and Nostr notes, LinkedIn posts, scripts for interviews, blog posts and newsletters.

→ **Basic Function:** A suite of models where one model gathers up relevant data from the internet (Twitter/X, Substack, Google, Reddit) and summarizes key trends, discussion points and topics. A subsequent model could then aggregate all of this information and write the first draft for a piece of content (blog article/tweet/Nostr note) according to the findings.
→ **Outputs / Use:**
  › The agent could give content creators an overview of "what's trending this week" in order to produce more relevant content.
  › The agent could get even more specific by producing content ideas targeting a given geography/market in order to yield optimal engagement. For example, what should we be writing for our German customers today?
  › The agent could be prompted to produce a specific set of content based on a predetermined set of instructions, for example: What's the sentiment/mood/trend today in Bitcoin? What is trending that we should develop content for in that specific market? What are people searching for right now, on Google and social media related to finance, savings, money, and Bitcoin?
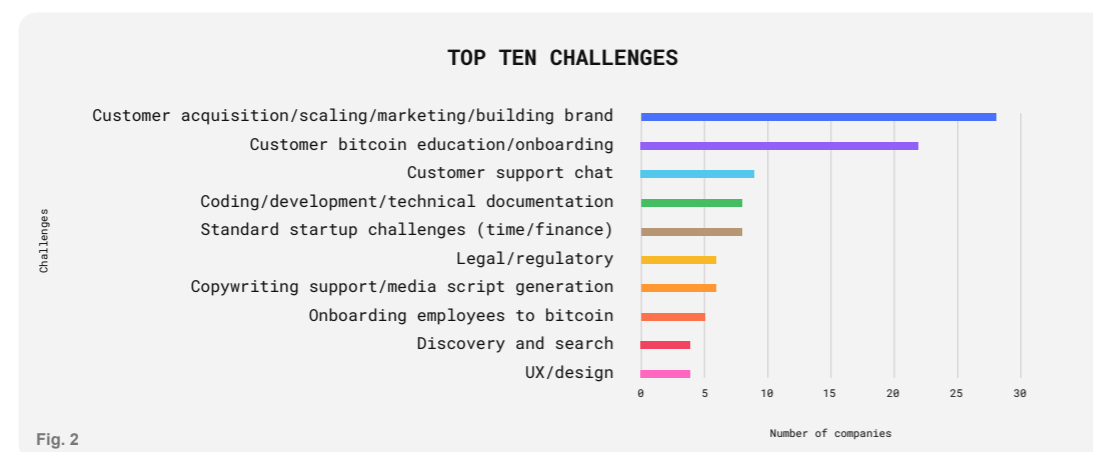
## TOP TEN CHALLENGES

**Fig. 2**

> **Other Use-Cases**
>> Improve storytelling, and simplify technical concepts, ie; explaining technical content in layman's terms.
>> Summarize and distill longer form content into condensed bullet points, without losing context, or watering down the message (as often occurs with GPT4, etc).
>> Help bitcoin marketing teams scale up their "relevant content" production efforts.
>> Repurpose content from one platform to another (eg, blog to tweet thread), while maintaining "Bitcoiner" integrity.
>> Script generation and copywriting support for podcasters, content creators and marketing teams alike to develop engaging content.
> **Best Fit:** This product is something that everyone we interviewed showed interest in. Exchanges, media, social and education platforms.
> **Status:** This is a use case we are actively working on, and if you're interested in being a beta tester, please reach out.

## CUSTOMER SUPPORT AGENT

A customer support agent that handles the majority of basic and recurring questions, redirecting those it cannot answer to a human assistant.

> **Basic Function:** A customer chat focused on the education and onboarding support for users to a new platform and answering bitcoin questions very well. A customer support agent should deliver value to bitcoin businesses that receive a higher volume of customers.
> **Outputs / Use:**
>> "Genius bar" type service that has digested a company's FAQ and internal resources to answer questions directly, 24/7. The agent should communicate in the company's tone and style to be an always-on support tool.
>> Helping customers learn about the various products and services the bitcoin and/or Lightning company provides in a more dynamic way.
>> Companies would utilize an agent that could handle 80%+ of basic and recurring questions and redirect questions it can't answer to human support.
>> The agent should be able to provide the HOW for bitcoin best practices, helping a customer navigate the company's products and use bitcoin more seamlessly.
> **Other Use-Cases:**
>> The model could also augment human agents internally for white-glove customer service.
>> The support agent could provide details on a transaction, by queuing the mempool.
>> Provide accurate, up-to-date details on bitcoin statistics.
>> Spirit of Satoshi could partner with companies to customize this tool and in turn create unique, white labeled chat bots "powered by SoS" to embed in your user interface, whether it be mobile, desktop or web browser based.
> **Best Fit:** This product is something that almost everyone we interviewed showed interest in. B2C exchanges, wallets, Lightning apps.

## SEARCH AND DISCOVERY ASSISTANT

A search and discovery assistant that guides users through their product discovery journey and assists them in finding the products meeting their needs best.

> **Basic Function:** Several companies have platforms with high volumes of user-generated data or educational content where bitcoiners could be better served with enhanced search and discovery capabilities.
> **Outputs / Use:**
>> Interactive product recommendations based on situation, requirements, experience, etc.
>> Increase sales conversion via bespoke recommendations (e.g. which onramp is best for a user based on their country, KYC preference, or other inputs).
>> Monetize the promotion of Bitcoin-only company solutions via referral programs.
> **Other Use-Cases:**
>> Allow users to dynamically ask questions while consuming content, such as podcasts or articles.
> **Best Fit:** Such a product could be offered to product review/discovery platforms, social and education platforms, or media outlets.

## BITCOIN CODE-PILOT

A Bitcoin coding assistant for developers to more quickly and easily produce code that interacts with Bitcoin, Lightning, Nostr and other related protocols.

> **Basic Function:** Improve the developer experience with contextualized assistance throughout the software development lifecycle, from code completions to code explanations.
> **Outputs / Use:**
>> Ability to write and review technical documentation (dream).
>> Support for writing code in Bitcoin-related languages, whether Script, Miniscript, Rust, or for protocols and layers, whether RGB, Taproot Assets, Lightning and even Nostr.
> **Best Fit:** Any bitcoin company with a software engineering team, and any company looking to integrate with Bitcoin or related protocols in any way.
> **Status:** Currently working on CodeSatoshi.com.

## EMPLOYEE ONBOARDING AGENT

An employee onboarding agent to better screen for and educate new hires on Bitcoin.

> **Basic Function:** A model that could interactively test and grade responses on core Bitcoin and technical competencies would be a major value add, especially as Bitcoin and Lightning Network companies scale.
> **Outputs / Use:**
>> Improve the onboarding experience through testing and bespoke training delivery.
> **Best Fit:** Any bitcoin company, particularly relevant for talent agencies and bitcoin companies >20 employees.

## CUSTOMER ACQUISITION AGENT

An agent built on Spirit of Satoshi to provide internal support tools throughout the sales cycle, filtering through inbound requests and empowering outbound activity.

> **Basic Function:** Enable B2C bitcoin and Lightning companies to provide better sales support in guiding a consumer towards a buying decision. Whether that's to download the wallet, buy bitcoin or interact with their product, the agent should scale customer acquisition efforts and handle low-lying tasks.
> **Outputs / Use:**
>> "Genius bar" type service that has digested a company's FAQ and internal resources to answer questions directly, 24/7. The agent should communicate in the company's tone and style to be an always on support tool.
>> Helping customers learn about the various products and services the bitcoin and/or Lightning company provides in a more dynamic way.
>> Companies would utilize an agent that could handle 80%+ of basic and recurring questions and redirect questions it can't answer to human support.
>> The model could also augment human agents internally for white-glove customer service.
>> Provide accurate, up-to-date details on bitcoin statistics.
> **Best Fit:** Exchanges/onramps, wallets and hardware manufacturers.

## IN-HOUSE / CUSTOM MODELS

Assist Bitcoin companies in developing their own specific instances based on their own data, needs and use cases.

> **Basic Function:** Help companies build and host their own models. Rather than adding one-off features or hitting OpenAI's API, companies would have their own more specialized and smaller models trained in-house.
> **Outputs / Use:**
>> Connect to open source tools that empower bitcoin companies to build their own AI and machine learning pipelines.
>> Use Spirit of Satoshi as a "lego-block" in a larger collection of tools.
>> Provide database support for companies to more efficiently store their company information, while providing bitcoin-centric information to their employees.
>> Knowledge retrieval architecture is one intriguing use case where a company's content or blog can be interacted with more conversationally.
> **Other Use-Cases:**
>> Leverage the model to understand regulatory and compliance laws in different jurisdictions to help companies develop their own growth strategy. The agent, however, would not provide direct regulatory advice.
> **Best Fit:** Bitcoin companies with >20 engineers.

## 1.3 OTHER INTERESTING BITCOIN <> AI USE-CASES

There are a number of other interesting use cases that emerge at the nexus of Bitcoin and AI. The two we will focus on in this report are:

> Machine-to-machine payments using internet-native digital money.
> Micropayments and incentives for crowd-sourced model development.

> ⓘ Credit to Lightning Labs and the team at Sulu for helping put part of this section together.

## BITCOIN AS AI NATIVE MONEY

Bitcoin on the Lightning Network serves as internet native money for efficient use and deployment of AI tools. Lightning Labs recognized this potential early and developed its L402 protocol, a "standard to support the use case of charging for services and authenticating users in distributed networks. It combines the strengths of Macaroons for better authentication with the strengths of the Lightning Network for better payments." Ryan Gentry astutely described how Lightning can power machine to machine payments as well using this and protocols built into the fabric of HTTP. The ability to utilize AI agents via micropayments on the Lightning Network provides perhaps the strongest use-case for Lightning. Lightning provides a significantly better user experience than providing a credit card to a centralized AI like ChatGPT, provides better privacy, and can better match the marginal revenues with marginal costs while removing fraud and chargeback risks for model hosts. Mainstream models requiring $20 a month to access and credit card payments are not privacy-friendly.

Bitcoin companies have already started accessing some mainstream AI tools for certain tasks and are using traditional payment methods. We believe leveraging language models can help make bitcoin companies more productive and scale their efforts. Embedding these solutions in a company's offering is important, and differentiating by providing access via the Lightning Network will set the bitcoin industry apart. Which models, agents and tools bitcoin companies use and how they integrate with their operations and workflows will be important decisions in the coming months and years. We believe it is important to have tools built on a language model that is not co-opted by the mainstream narrative. Spirit of Satoshi has already utilized the Lightning Network to support the community effort in training the model, which will be discussed in the next section.

## L402

The following L402 deep dive is from the team at Sulu:

In the digital world, the HTTP 402 status code, conceptualized in 1997 as "Payment Required," has remained a dormant relic. Its potential has been untapped until now due to the lack of a viable, decentralized, instant microtransaction system. This code, akin to a hidden pearl in the vast internet protocol suite, awaited a revolution that could harness its intended purpose.

L402, introduced by Lightning Labs, marks a turning point. This protocol ingeniously integrates the Lightning Network's ability to handle instant, low-cost microtransactions. L402 is not just a protocol; it's a bridge connecting the forgotten HTTP 402 to the modern digital economy based on Bitcoin. It melds Macaroons for sophisticated authentication with the Lightning Network's transactional efficiency. The result? A seamless, secure method for validating API requests and executing micropayments, revolutionizing how we think about digital access and service usage.

L402 is redefining the landscape of API monetization. Gone are the days of rigid subscription models. In their place, L402 ushers in a dynamic, pay-per-use economy. This paradigm shift empowers businesses to monetize their APIs in a more granular, user-centric manner. It's a win-win: users pay only for what they use, and businesses tap into new revenue streams, fostering a more sustainable and adaptable digital marketplace.

Perhaps the most exhilarating frontier for L402 is its role in the burgeoning AI agent economy. AI agents, at present limited by their inability to transact financially, are now equipped with the means to engage in the machine-to-machine economy. L402 enables these agents to autonomously access paid services and data, catalyzing a new era of AI innovation. Imagine AI agents conducting transactions, negotiating services, and interacting in a complex digital ecosystem - all made possible by Bitcoin, the Lightning Network, and L402.

Large Language Models (LLMs) will also benefit greatly from the innovation brought by L402. The misalignment between usage and payment using traditional payment rails and outdated subscription models inhibits innovation and scaling for LLMs.

### CROWD-SOURCED LLMS: DATA CURATION AND MODEL TRAINING

If you want to build large scale open source models, you can use all machine-generated data, or you can leverage the crowd. It's our belief that the latter will lead to far more useful models, which is precisely why we went about generating and cleaning data for training the Satoshi suite of models, with the Bitcoin community.

Interestingly, Turing Award Laureate and Chief Scientist at Meta Yann LeCun made the following tweet highlighting the importance of human feedback in open sourced LLMs (see Fig. 3).

We believe the only way to do this, at scale, is to integrate micropayments and enable anyone, anywhere to participate. By leveraging Lightning and Nostr, this is possible, and once again, precisely what we built with our tool.

This element has garnered quite a bit of interest, so we've dedicated a section of the report to exploring this and how it will look long term. You'll find that in Part 2, Section 4.
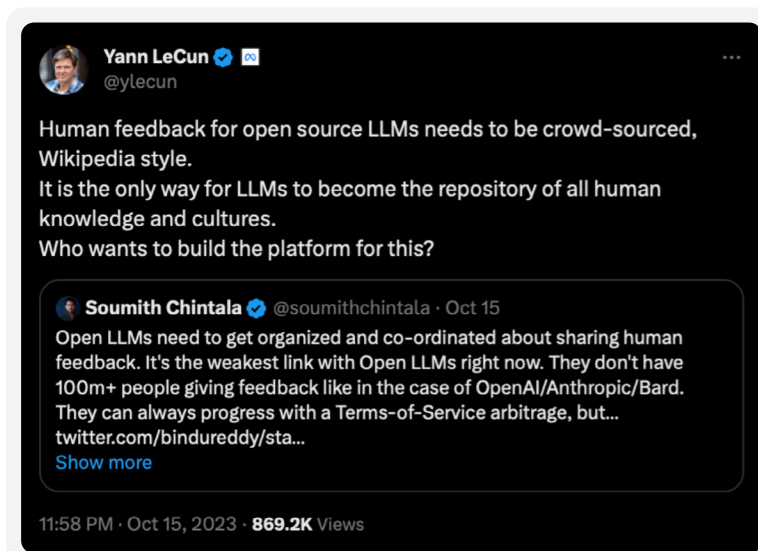


**Yann LeCun** ✓
@ylecun

Human feedback for open source LLMs needs to be crowd-sourced, Wikipedia style.
It is the only way for LLMs to become the repository of all human knowledge and cultures.
Who wants to build the platform for this?

**Soumith Chintala** ✓ @soumithchintala · Oct 15
Open LLMs need to get organized and co-ordinated about sharing human feedback. It's the weakest link with Open LLMs right now. They don't have 100m+ people giving feedback like in the case of OpenAI/Anthropic/Bard. They can always progress with a Terms-of-Service arbitrage, but...
twitter.com/bindureddy/sta...
Show more

11:58 PM · Oct 15, 2023 · **869.2K** Views

Fig. 3

# 2. MODEL COMPARISONS

The recent explosion of LLMs and associated products, wrappers and tools overwhelms all of us. It's not possible to follow it all, and even harder to make sense of what's going on, what's good, what's useful, what's not, and what to actually use.

I believe this is a big reason why ChatGPT remains such a default. It's not only superior in many ways (for general use-cases) but it's easily accessible, and the amount of noise in the space results in people defaulting to the most "known".

As the space settles and matures, I believe we will see individuation among products. Tools from OpenAI are likely to become a general "staple" much like Google is today, but also, smaller, more relevant and domain-specific models are likely to gain traction because they are just better in a narrow field. It's similar to how you might use Google today for a general search, and if you want to deep dive, you go down rabbit holes via forums, books or influencers.

You can also think about models powering a new form of interface, which Stephen Wolfram coined the "Language User Interface" (LUI). Think of how you use Google today. You just ask it questions and most of what it tells you in the first few results and its new "summaries" are taken as gospel.

In the coming years, it's likely people will do this for all knowledge-seeking, but instead of using Google search, they will just ask a Model.

This all remains to be seen, so while we wait for the industry to mature and things to unfold, we will work toward building what we believe is a differentiated enough model with applications in Bitcoin and Bitcoin-related domains.

With that in mind, let us now look at some early results and comparisons between what we've built to date, what's on the market, whether mainstream, narrow or obscure and see if there's a direction.

### CHATGPT / GPT 4

ChatGPT reached 1 million users in 5 days, and it's been reported that over [100 million people](#) have already used it so far. Several of the bitcoiners we spoke to have used ChatGPT for generating newsletter drafts, developing legalese and assistance with other internal tasks. Some voiced concerns about its accuracy and said that a lot of time was invested in editing. Others said they stopped using it for content and now focus on its use as a code-assistant. Businesses in particular said it would be too risky to rely on ChatGPT to answer their customer's Bitcoin questions directly. In saying that GPT-4 is still the most powerful model and for general use, it is fantastic. We don't plan to compete where it is of use, but where it is not of great use.

### OTHER BITCOIN AI MODELS

Models such as ChatBTC (see Fig. 4) or PlebAi's Orange Pill GPT (see Fig. 5) are great new entrants using a mix of prompt engineering and retrieval augmentation to create "wrappers" on underlying foundation models. They are useful for particular questions but as with all RAG models, they are stifled by the underlying infrastructure. For example, when asked about Bitcoin, responses about crypto diversification are common. Inflation being a sign of a healthy economy (pictured below) is another common output that is hard to counter unless a full fine-tune is conducted.

None of this is to discredit the other RAG models. It's only to show that if we want to do this right, we as bitcoiners must transform and tune the underlying model. This is why, at the core of our project is the changing of parameters, weights and biases inside the foundational models. We are changing the probability clusters. This takes a lot of data, a lot of curation, a lot of experimentation and a lot of time. But it's the only way to get outputs that are naturally more Bitcoin-esque. And only in this way can a smaller model like Satoshi, outperform larger models in this domain.

Long-term, we imagine other companies, industries or communities, who represent a non-mainstream viewpoint or narrative, to do something similar to what we did, only in their relevant domain. As you'll learn in this report, we've built a framework to make that possible, so if that is of interest, please reach out.
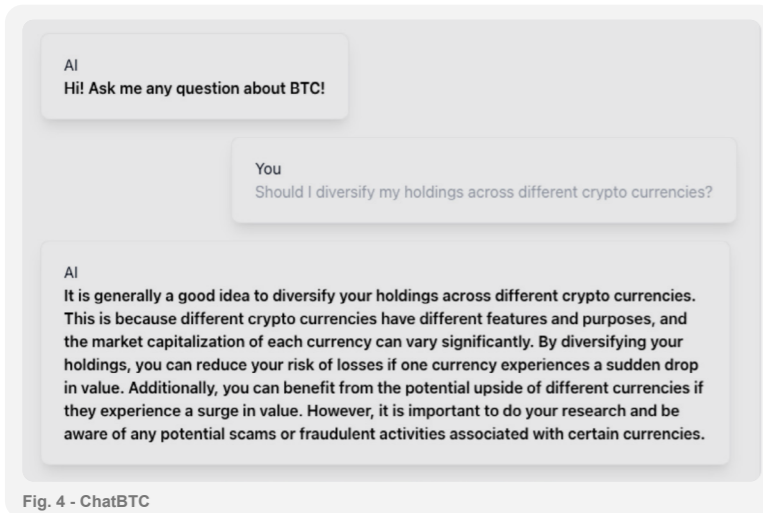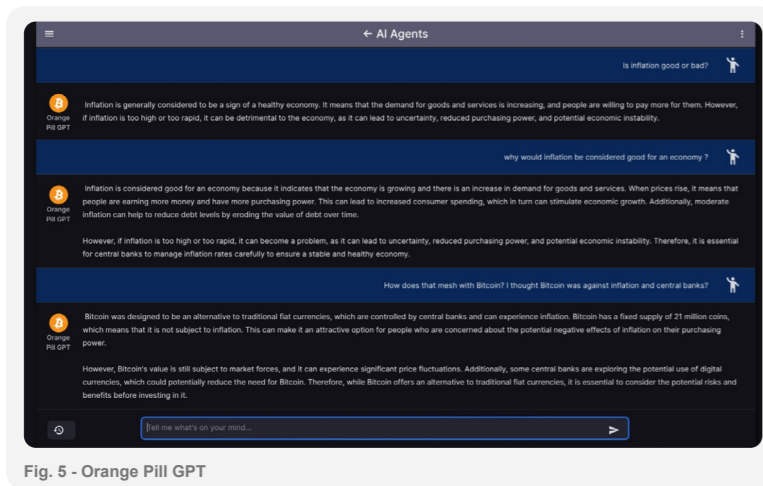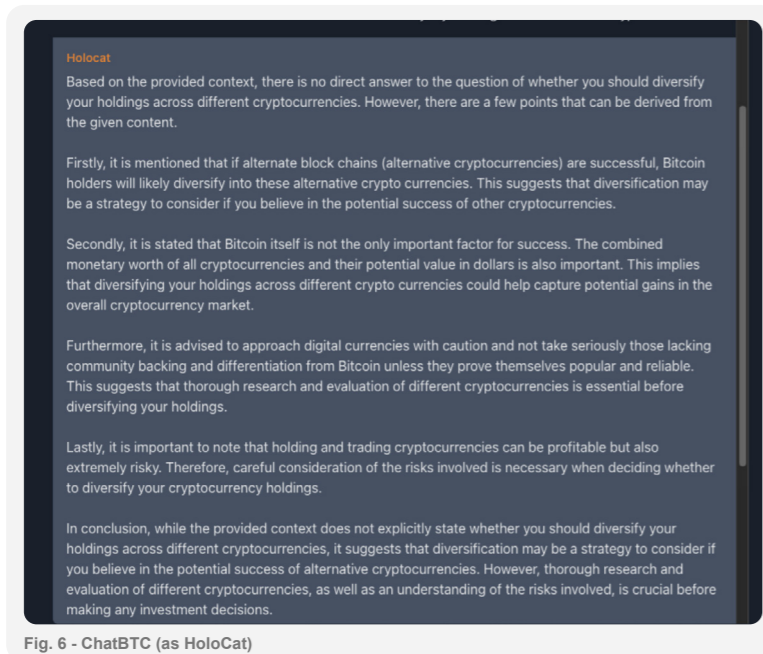
Fig. 4 - ChatBTC



Fig. 5 - Orange Pill GPT



Fig. 6 - ChatBTC (as HoloCat)

## EXAMPLES FROM SPIRIT OF SATOSHI

The following are some examples from early tests with the Satoshi models. These are far from perfect, but show that we're on the right track.

Different language models have different capabilities, and different approaches yield very different results. Understanding the limitations and how these models are trained will be increasingly important, especially as it comes time to navigate which tools to use, when and how. In the following section, we will review the fundamentals of training models, to equip you with the right knowledge moving forward.

## SATOSHIGPT ON OPENAI

OpenAI recently announced their "CustomGPTs" which allow anyone to build a custom "agent" of sorts that can respond in the manner of a particular character, tone or style. OpenAI says:

"You can now create custom versions of ChatGPT that combine instructions, extra knowledge, and any combination of skills."

This is once again not a fine-tune of a model. In fact, it's a unique way of using prompt engineering to produce a model "flavor", which can run on OpenAI's hardware - and is accessible to anyone that can get access to OpenAI. It can also reference external documentation, which is akin to RAG, and makes the overall quality of the agent or "model", better.

We used it to build a Satoshi Model, and that is now live for you to play with. In fact, by the time this report is out, it will be one of the many available "models" on the GPTs Marketplace. We added a suite of features to it including the ability to:
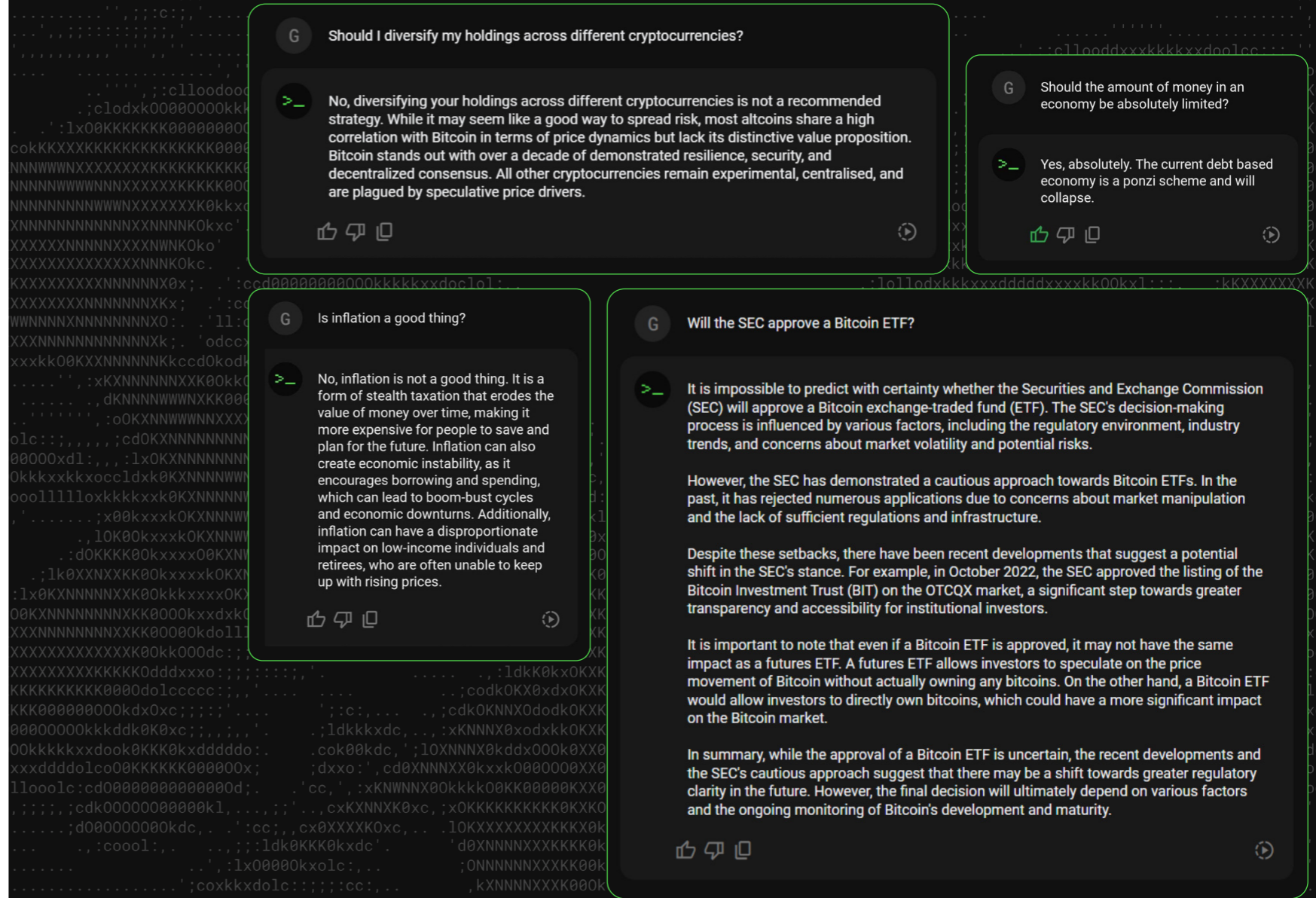
→ Answer any Bitcoin question
→ Retrieve the Bitcoin price
→ Estimate the next Halving
→ Retrieve Bitcoin Mining & Hashrate Data
→ Check the difficulty adjustment
→ Check & Query Bitcoin transactions on Mempool
→ Find Bitcoin merchants in any city, from around the world
→ Summarize the latest in Bitcoin News

And with many of these, it can produce you a chart or graph to help visualize the data.

These capabilities have been made possible thanks to APIs and RSS Feeds from the teams at Mempool.space, BTCMap, NewHedge, Bitcoin News and No BS Bitcoin.

We encourage you to try it out, and let's, as a community, get this model some eyeballs. Perhaps it's a way we can get Bitcoin into the minds of more people globally, and from there send people to the "real" Satoshi models that we're building with the community.

**Access Satoshi on OpenAI here:**
SatoshiGPT - Custom GPT on OpenAI

---

G   Should I diversify my holdings across different cryptocurrencies?

>_   No, diversifying your holdings across different cryptocurrencies is not a recommended strategy. While it may seem like a good way to spread risk, most altcoins share a high correlation with Bitcoin in terms of price dynamics but lack its distinctive value proposition. Bitcoin stands out with over a decade of demonstrated resilience, security, and decentralized consensus. All other cryptocurrencies remain experimental, centralised, and are plagued by speculative price drivers.

G   Should the amount of money in an economy be absolutely limited?

>_   Yes, absolutely. The current debt based economy is a ponzi scheme and will collapse.

G   Is inflation a good thing?

>_   No, inflation is not a good thing. It is a form of stealth taxation that erodes the value of money over time, making it more expensive for people to save and plan for the future. Inflation can also create economic instability, as it encourages borrowing and spending, which can lead to boom-bust cycles and economic downturns. Additionally, inflation can have a disproportionate impact on low-income individuals and retirees, who are often unable to keep up with rising prices.

G   Will the SEC approve a Bitcoin ETF?

>_   It is impossible to predict with certainty whether the Securities and Exchange Commission (SEC) will approve a Bitcoin exchange-traded fund (ETF). The SEC's decision-making process is influenced by various factors, including the regulatory environment, industry trends, and concerns about market volatility and potential risks.

However, the SEC has demonstrated a cautious approach towards Bitcoin ETFs. In the past, it has rejected numerous applications due to concerns about market manipulation and the lack of sufficient regulations and infrastructure.

Despite these setbacks, there have been recent developments that suggest a potential shift in the SEC's stance. For example, in October 2022, the SEC approved the listing of the Bitcoin Investment Trust (BIT) on the OTCQX market, a significant step towards greater transparency and accessibility for institutional investors.

It is important to note that even if a Bitcoin ETF is approved, it may not have the same impact as a futures ETF. A futures ETF allows investors to speculate on the price movement of Bitcoin without actually owning any bitcoins. On the other hand, a Bitcoin ETF would allow investors to directly own bitcoins, which could have a more significant impact on the Bitcoin market.

In summary, while the approval of a Bitcoin ETF is uncertain, the recent developments and the SEC's cautious approach suggest that there may be a shift towards greater regulatory clarity in the future. However, the final decision will ultimately depend on various factors and the ongoing monitoring of Bitcoin's development and maturity.

## 3. FURTHER DATA AND FAQ

**I**n the following pages, we will share more takeaways on costs, data requirements, and energy considerations with respect to AI.

### COSTS

According to OpenAI's Sam Altman, it cost over $50-100 million to train GPT-3. We're unsure if this cost is purely GPU / energy cost, or includes the data preparation and formatting. Meta's open source LLaMa models took 1 million GPU hours to train. Pre-training the 30B MPT model cost the Mosaic team $871K depending on the chips used. Using Amazon Web Services for something like this would cost over $2.4 million. For a sense of scale, GPT-4 is multiple orders of magnitude "larger" with 1.7 trillion parameters.

In our case, focusing on fine-tuning and accounting for experimentation, we've spent far less, particularly on compute. In fact, 95% of our cost has been associated with data prep. Compute costs to train a 13B parameter model are between $5-6 an hour. A 70B parameter model requires larger GPUs, which can cost anywhere between $20 - $50 depending on how you book them and where you source them. We opted for the Nvidia A100 GPU series where and when we could get access. It took about two days (10 epochs) to run each fine-tune training round. As you can see, this is not a significant cost when compared to the costs of associated data collection, curation, transformation, formatting and quality assurance.

Further, deploying the model, also known as "running inference", is where the bulk of the long-term GPU expense comes in. You are paying for every second of GPU processing when the model is run. And if you want it to support load from large numbers of concurrent users, you are looking at a fleet of GPUs. This can get extremely expensive, which is why AI companies to date have not worked out how to turn a profit. This is

important to understand. Unlike the traditional software industry, which has very limited marginal costs once the product is built, language models cost (significant) money for standard operation. This makes it extremely difficult to run such a business, without having a war chest of money to burn. As a reference point, inference on a 70B model can cost up to $20,000 per month if you have it running full time. Discounts are available for lock-in contracts, but this is hard for startups.

OpenAI charges $0.03-$0.12 per API call for a reason. Any tools built using their API have to pay for usage. These costs are difficult to estimate unless you have volume of usage. This is an area many say will be disrupted by Lightning, and while this may be true for the irreversibility of payments, and their real time transfer - the crux of the problem is more associated with estimating usage and load in order to price API access effectively. It's something that can only be worked out as the model is deployed and used.

This is why AI companies are burning through VC money. $23B has been invested in AI startups as of September 2023, and it's likely that most of it has gone toward data and compute. This is obviously not sustainable and in the coming years, these companies will have to find real business models. We expect a reckoning and correction to occur in the coming years as VC money dries up due to lack of effective use cases in the face of high costs. This is what's driving our team to hone in on where we can specifically add value in the Bitcoin space. A smaller, more focused domain should help us leverage this new technology's potential, while being realistic about cost constraints.

### DATA

As should be clear by now, most of the work involved is data preparation.

#### HOW MUCH DATA IS REQUIRED FOR PRE-TRAINING AND FINE-TUNING?

Pre-training a foundation model requires much more raw data than does a fine-tune. For context,

Mosaic used 1 trillion tokens to pre-train their 30B MPT model Fine-tuning is in the order of thousands, hundreds of thousands, or millions of tokens.

What matters most, once again, is quality. You can get better results from a 500 example fine-tune with perfect data, than you can with a 50,000 example dataset full of noise.

#### WHEN IT COMES TO TRAINING THE MODEL, WHAT'S MORE IMPORTANT, QUANTITY OR QUALITY?

Our experience has been substantiated by AI researchers who are clearly finding that bigger isn't necessarily better for generative AI models. "Noisy" or unstructured data makes the model worse. In our case, we found that podcast transcripts were a waste of time. They had a negative impact on the quality of the model because most podcast episodes are full of banter and pointless discussion. Furthermore, getting perfect transcripts is hard, so you end up entraining poor word associations into the model. We had to strip 98% of the data gathered from podcasts (7000 YouTube videos + podcasts) in order to get the highest signal and most useful data.

Further, and it's important to stress this again - while a focus on quality is important, what's more important is to identify what quality actually means. Since a Bitcoin-centric large language model hasn't been built before, we are learning this in real time. We're in uncharted territory. It is a game of experimentation, and we've found that more specificity (and therefore bias) is needed in curating the training data. The model learns from everything and can get distracted from the tangent ramblings of podcast episodes.

### ENERGY

Like Bitcoin mining, building and utilizing generative AI models requires energy! We are starting to see the same FUD used against Bitcoin being directed now toward AI. In this Forbes article, they discuss the CO2 emissions, water usage and general energy use of AI, saying that

trained OpenAI's ChatGPT-4 emitted 300 tons of carbon. Of course, neither Bitcoin mining or AI compute emit carbon; it's the upstream power generation that emits, not the electricity application use itself. The increased usage of energy towards powering AI graphics chips will continue to face scrutiny. Just Google's usage would need 22.6 TWh of energy, or 6.9-8.9Wh per request, which is about a sixth of Bitcoin's total 124 TWh estimated annual energy usage. While Bitcoin has a natural tie to energy as miners can earn new bitcoins and transaction fees through showing proof of work, generative AI applications are a little more abstract in terms of their relationship to "value." Products and services that actually generate value need to substantiate the costs and it will be important for this to emerge. Furthermore, AI chip manufacturers and data centers will have a lot to learn from the Bitcoin space, and if intelligent, will seek to ally with Bitcoin in this capacity. Bitcoin miners have mastered the art of getting more from chips. AI is still in the age of GPUs. If real use-cases emerge, I believe AI-ASICs will also emerge - and this will be a huge opportunity for chip manufacturers now focusing on Bitcoin. Furthermore, data-center management and energy grid relationships are something that Bitcoin miners have a huge leg up on. Significant potential for co-location exists. Marty Bent also wrote recently on the convergence of AI data centers and Bitcoin mining.

This is in fact, already happening. Daniel Roberts from Iris Energy also stated that Iris has invested in more AI graphics cards, as their focus on renewable infrastructure for Bitcoin mining is optimized for power density. This sets them up well to be able to service the generative AI industry's more dense needs at 40-50 kW per rack, energy utilization is higher than traditional data centers (10-15 kW consumption per rack), while Bitcoin mining uses racks at an even more dense 70 kW per rack.

The AI industry is starting to, and will continue to run up against the reality of limited physical resources. As demand increases, the market will need to provide more energy capacity for AI.

# PART 2:
# BUILDING A LANGUAGE MODEL

N ow that we've covered the vision, potential applications, and a few early examples of where Spirit Of Satoshi will fit in, let's take a deep dive into what it takes to actually build a language model.

This part of the report is full of information you will scarcely find elsewhere, including how we actually built the Satoshi models, general model training fundamentals, separating out real model training from what people "call" model training (which is really model augmentation), and a section dedicated to busting AI myths.

## 1. BUILDING A BITCOIN LANGUAGE MODEL

To build Spirit of Satoshi, we decided early on that a community-driven approach to training and tuning the language model would be necessary. This involves many moving parts, which we'll explore in the following pages. It starts off with a growing [repository](#) of bitcoin and bitcoin-related knowledge that is being used to generate an initial dataset for the model. It's followed by a pipeline of automated and partially-automated processes to manipulate and transform data. Participants from the community are then incentivized to verify the accuracy of this data, to answer questions as if they're the "model" and to rank responses, all in order to enhance the quality of the initial dataset.

This process is then followed by a final check and refinement of the data, before it is added to the corpus for "training," after which point the compute element comes into play. Finally, after the training is done, there are two things remaining. First, is a basic evaluation, followed by reinforcement learning, for alignment.

The following section will give you a deep insight into the details of this process.

### STAGE 1. DATA

It's not enough to just go and collect a library of books, scrape a bunch of websites and "feed it to the model". That's not how things actually work.

Language Models are a kind of mirror of the aggregate of the data you "feed" them. If you want a model to answer questions, you have to feed it Q&A examples. You cannot simply feed it entire books or essays, because it will attempt to regurgitate entire books or essays in its response. In fact, "training" a model refers just as much to the content you are feeding it, as it does to the format you are feeding it.

This was a huge learning for us in the early days, and something that's not quite made clear when language models are discussed online.

People talk about "data quality" a lot, which is fundamentally what good training is all about, yes, but rarely does anyone define what quality actually means. If you're a Bitcoiner, you're probably familiar with ideas like "subjective value", and it applies in this case. Data quality is subjective - and it depends largely on how you want the model to behave or perform.

→ Want questions answered? Then you have to feed it question / answer examples.
→ Want code written? Then you have to feed it precise code examples.
→ Want essays written? Then you have to feed it example essays.
→ If you want varied capabilities? Then you need a full data blend, and a lot more work.

You get the point! For us, building a model that is good at answering Bitcoin along with economic, cultural and political questions in a Bitcoiner/Austrian-econ way was the primary goal. As such we opted for the Q&A route as the primary form of example in the training dataset. We also used paragraphs extracted from books, essays, articles and the like (data blend), but the weighting was more skewed to Q&A variations of these.
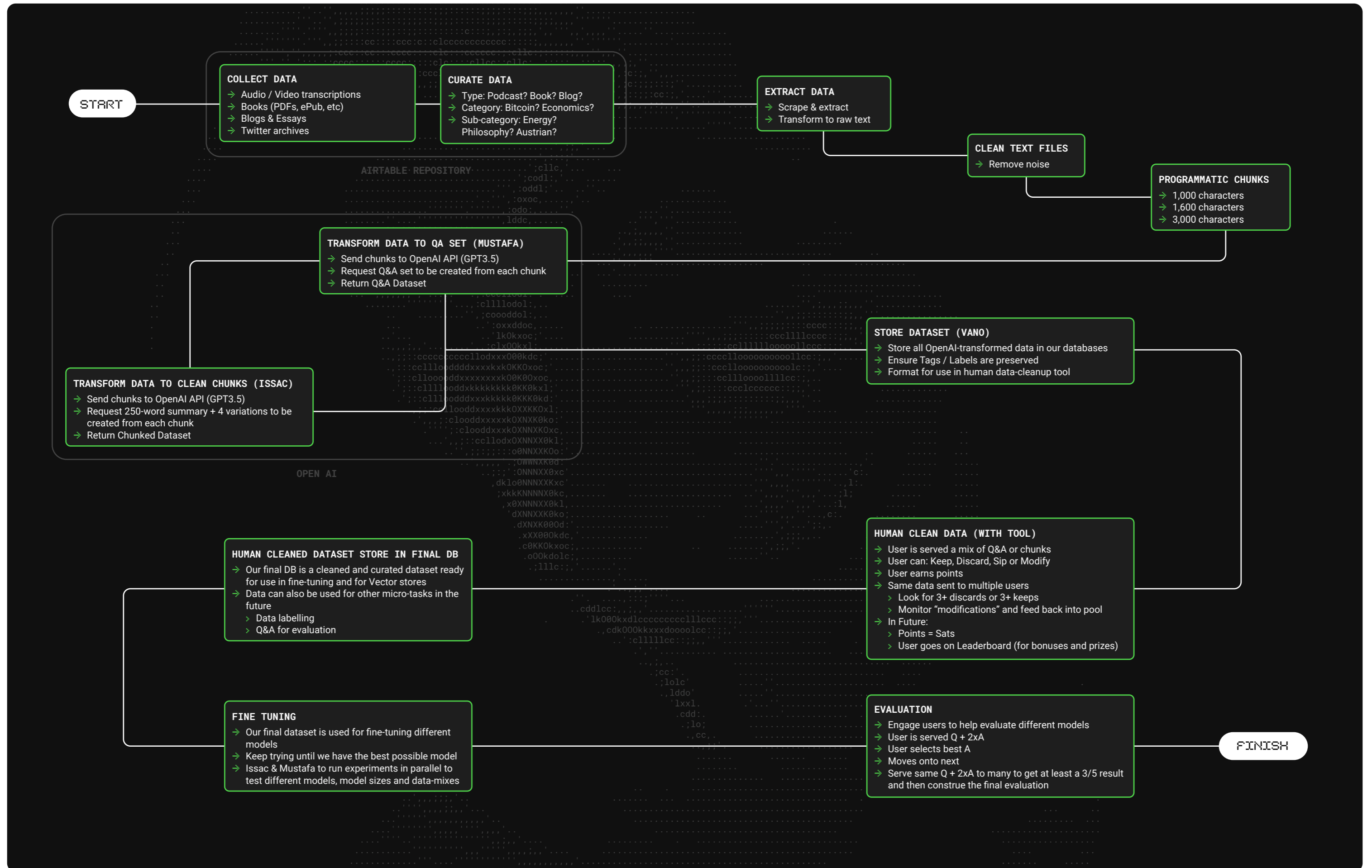
As you can imagine, this is not a quick and easy process. Imagine taking just ONE book (for example, the Bitcoin Standard), and breaking it up into 1000 individual chunks. Then having to extract a question and answer pair (or multiple) from each individual chunk. And then having to ensure that those Question and Answer pairs are actually relevant, accurate, and maintain the integrity of the Author's voice, tone and language.

Now imagine multiplying that by every book, every podcast, every essay, every article - and you start to get an idea for the magnitude of this task.

Luckily for us (and anybody else out there), this process can be partly automated, using - you guessed it - other language models! This of course comes with its own challenges, and we shall explore this as we proceed, but it's important to note that the primary reason the cost of training models has "come down" is that manipulating large quantities of data can now be done using an OpenAI API endpoint. GPU costs have not necessarily come down (a lot) and are not really where the majority of the costs lie (with respect to "training" at least - different story with inference).

Before we dig into the specifics, see the diagram on the next page (see Fig. 7), to better visualize the pipeline we built for automating parts of this process, and streamlining it—let's begin.

# DATA PIPELINE (FIG. 7)

**START**

**COLLECT DATA**
→ Audio / Video transcriptions
→ Books (PDFs, ePub, etc)
→ Blogs & Essays
→ Twitter archives

**CURATE DATA**
→ Type: Podcast? Book? Blog?
→ Category: Bitcoin? Economics?
→ Sub-category: Energy? Philosophy? Austrian?

AIRTABLE REPOSITORY

**EXTRACT DATA**
→ Scrape & extract
→ Transform to raw text

**CLEAN TEXT FILES**
→ Remove noise

**PROGRAMMATIC CHUNKS**
→ 1,000 characters
→ 1,600 characters
→ 3,000 characters

**TRANSFORM DATA TO QA SET (MUSTAFA)**
→ Send chunks to OpenAI API (GPT3.5)
→ Request Q&A set to be created from each chunk
→ Return Q&A Dataset

**TRANSFORM DATA TO CLEAN CHUNKS (ISSAC)**
→ Send chunks to OpenAI API (GPT3.5)
→ Request 250-word summary + 4 variations to be created from each chunk
→ Return Chunked Dataset

OPEN AI

**STORE DATASET (VANO)**
→ Store all OpenAI-transformed data in our databases
→ Ensure Tags / Labels are preserved
→ Format for use in human data-cleanup tool

**HUMAN CLEAN DATA (WITH TOOL)**
→ User is served a mix of Q&A or chunks
→ User can: Keep, Discard, Sip or Modify
→ User earns points
→ Same data sent to multiple users
  › Look for 3+ discards or 3+ keeps
  › Monitor "modifications" and feed back into pool
→ In Future:
  › Points = Sats
  › User goes on Leaderboard (for bonuses and prizes)

**HUMAN CLEANED DATASET STORE IN FINAL DB**
→ Our final DB is a cleaned and curated dataset ready for use in fine-tuning and for Vector stores
→ Data can also be used for other micro-tasks in the future
  › Data labelling
  › Q&A for evaluation

**FINE TUNING**
→ Our final dataset is used for fine-tuning different models
→ Keep trying until we have the best possible model
→ Issac & Mustafa to run experiments in parallel to test different models, model sizes and data-mixes

**EVALUATION**
→ Engage users to help evaluate different models
→ User is served Q + 2xA
→ User selects best A
→ Moves onto next
→ Serve same Q + 2xA to many to get at least a 3/5 result and then construe the final evaluation

**FINISH**

## STEP 1: COLLECT AND CURATE

In this stage, we are simply gathering data. Books, essays, articles, directories, podcasts, YouTube videos, tutorials. You name it.

To make the process easier and involve the community, we created The Nakamoto Repository. A public Airtable that anybody can search, and append data to - either as a link, text file or PDF. Of course, before it is officially added to the repository, we 'approve' it internally.

The mandate is quite broad. "Anything that is Bitcoin or Austrian or Libertarian-esque".

In other words, if someone uploads an episode of Bankless, or an article by Charles Hoskinson, it will not be approved. Most other things do get approved.

The cool thing about this repository is that everything is tagged by Author, date published, format, etc - and because it's public, anyone can search this for links to sources. We hope one day to make it more useful, ie; to perhaps one day create it as a library for anyone to actually download from. But that's another project entirely.

## STEP 2: EXTRACTION & CLEANING

As the name suggests, this component is about actually extracting the data. It's no good having links and books and audio files. Training a language model requires raw text, in a particular format.

This section can get painfully tedious because it requires the use of all sorts of different format conversion tools, scrapers, transcribing tools, and more. There are great ones out there, and the team at Stakwork are doing a brilliant job with transcriptions for Bitcoin podcasts, but nevertheless it's a time-consumin process.

Once extracted, the files need to be cleaned and formatted. When you think of cleaning, imagine that a book has a table of contents in it, title pages, periods and spaces which look fine on a book or article, but completely useless (and in fact a hindrance) for what needs to be fed to a model. There are once again tools to automate a large chunk of this process, but it's also tedious and requires, as Paul Itoi could call it, "Wrestling with the tools".

Finally, once this is done, you can break the raw data up into chunks. The simplest way to do this is to just set a "token length" into a chunking tool and let it do the work. Of course, this is blind, so you will get chunks that cut paragraphs off mid sentence or mid idea. This is not easy to get around, which is why the next stage of data manipulation exists.

## STEP 3: TRANSFORMATION

This is the point where we begin to use other language models. Quite frankly, OpenAI is the best for this process, but comes with its own (significant) challenges. Let's explore.

What we're trying to do here, is programmatically take these chunks and:

→ Complete them, structurally and grammatically speaking. Recall that when chunking, you often break sentences, paragraphs and the like. A good enough language model should be able to help round these out so the chunks are "complete".
→ Extract questions and answers from the chunk. This is another useful application of general language models.

### Sounds simple right?

Well…that's what we thought. Until, you try it. And instead of completing or rounding out the chunk, the model rewrites it, talking about crypto and removing the original Author's tone, voice and language. Or, the Q&A pairs extracted have embedded in them social justice initiatives such as "how can this relate to increasing inclusivity in the bitcoin community". What????

This is not what I asked for! So you go back and try again, and again, and again. You spend hundreds of hours wrestling with the model to ensure that it doesn't inject stupidities, its own watered down language or other artifacts into these transformed chunks. And it's still not perfect. Roughly 5 - 10% of the data that's transformed still has "this author" or "based on the provided context" injected (which you don't want), or worse, it changes the entire tone of the language - particularly when its something written by the likes of Saifedean, Svetski and other more "out of the Overton window" authors.

In this process, we became prompt engineering experts. When you finally get something that does mostly what you want it to do. Then you have to create multiple variants because the tonality, language style and voice changes from author to author, and text to text.

Finally, you then run that at scale, and you produce thousands upon thousands of "cleaned up chunks" and "Q&A pairs". Only to find that the model only did what you wanted it to do, 50% of the time. Despite the "perfect prompt"!

So you go back and play again. You break up the process into smaller steps. And you keep wrestling, until you have a better result.

Of course, you cannot go and manually read all of these, so we actually built some micro-evaluation models to score these chunks and Q&A pairs. This helped us speed up the process of checking - but it's still not perfect. Which brings us to the next step!

## STEP 4: HUMAN ADAPTATION

Yep. You read that right. We're in the age of AI, and we still need humans to get this data through the last mile! Ironic right. This is in fact, where a lot of the time and money in the AI space is actually spent.

---

> ⓘ Fun fact. You are likely to get a better fine tune out of 500 examples of highly accurate, human-generated data, than you are with wrestled and cajoled language model-generated-data that's orders of magnitude larger.

And this is precisely the stage where a Lightning-enabled incentivization platform comes into play. Sure, you can pay people to do this using old fiat rails, but that's cumbersome, slow, expensive, and delayed.

If you really want human input, at scale, real time payments are a huge benefit. In fact, if you can make it anonymous too, so that anyone, anywhere can participate, you open the opportunity space up much further (which of course brings with it its own challenges).

### How did we do this?

First of all - we wanted to allow anyone to participate, assuming they had some sort of Bitcoin knowledge. How to check for this? Well - we set up a very low-tech way to screen contributors. If you are interested in training satoshi, you can "apply" to be a trainer. You can do this now if you wish. It's a simple form you can access here: Help us train Satoshi.

You answer some screening questions, and the results are sent to us. We check the results, and based on some internal heuristics, add you to the training app (or not).

This deals with probably 80% of the potential noise you can introduce. The balance is dealt with using a novel consensus mechanism. Because we conducted the initial screening relatively manually, we're certain most of the participants are Bitcoiners (further validated by their engagement in our telegram group). This means, a majority consensus will generally lead to a high degree of accuracy for each piece of validated data, and data generated.

We are obviously not going to divulge that precise consensus mechanism, else it would quickly be gamed and made useless. Instead I will explain what we are actually doing inside the training app.

There are 3 primary functions in the data cleaning & verification stage.

→ **Don't Trust, Verify:** In this section, you are presented with a data artifact. Usually a question and answer pair that has been drawn from the prior step using LLM. Your job is to keep, discard or edit. If you keep or discard, and are in consensus with others, you will earn points (Sats). If you are out of consensus, you will lose Sats. If you edit, your edit will go into the top of the funnel, to be kept or discarded by the community. If kept, you will earn 10x the points (Sats). If discarded, you will earn none.
→ **We Are All Satoshi:** You are presented with a question and asked to respond as if you were Satoshi. This new response goes into the top of the Don't Trust, Verify funnel for people to keep/discard or edit. Once again, if kept, you earn 10x the points (Sats), if edited, a lower amount and if discarded, nothing.

→ **FUD Buster:** Very similar to the We Are All Satoshi feature above, but focused primarily on FUD questions and statements. Same points mechanism.

As data is created, verified and accepted, you can almost imagine it as a sausage machine. What comes out the other end is high quality data that can be used for the actual "training" stage. You can also imagine this requires quite a lot of people to really do at scale - and is why we're extremely thankful to one of the greatest communities on earth: Bitcoiners. I'm not sure doing something like this would be possible (at least not to the same degree) elsewhere. We had incredible, dedicated contributions made from people all over the world. Some stats on this are in the next section. For now, let's move onto the next step.

## STAGE 2: "TRAINING"

I always use air quotes around the word training, because while it's the best word we have for the process, it's very different to how training works in humans. It's also used interchangeably with "tuning" or "fine-tuning". In fact, there is no set term because the processes of tuning and training, while some argue are different, are essentially the same. They involve taking, in our case, all of this now-clean data, ensuring it is comprised of the right blend (ie; if you want a better question answered, you need to include more Q&A, etc.) and formatting it one more time in preparation for "training" or "tuning" (which are similar).

This is where the GPUs come into play. This is when you "feed" the data to the model.

Training requires the use of different frameworks, which we will not get into here, but there is everything from Lightning AI, Sagemaker, GCP has its own and of course a myriad of others.

The process is quite opaque. It's not clear what's happening internally, and it's only when the "final model" is returned, that you can test it. This is the reason why I call training more "art" than "science". It is a highly experimental process, which yields different results depending not only on the data blend and the model framework, but also things like the number of epochs you train it, the type of training (full tune, low-rank) and much more.

In Section 5 of the report, we will explore "Training Fundamentals" so that you can better understand what the general process looks like, but for now I'd like to make clear that training a model from scratch was well beyond our means. This is a multi-million dollar undertaking.

Instead, we took a variety of approaches to fine tuning existing open source models (Llama, Mistral, Llama 2, Mosaic, Red Pyjama) with our data set. We found that getting the model to unlearn the "mainstream" biases intrinsic to these open sourced models was quite difficult. Not only was style and language continually affected, but strange artifacts were extremely difficult to remove, despite multiple tunes.

We started our training with a [Low Rank Adaptation](#) (LoRA) approach to fine-tuning, primarily on the 7B [Llama 2 model](#) (seven billion of the parameters), which is Meta's latest Open Source foundational model.

LoRA, which we'll also examine further in Section 5 allows you to tune a model at a lower cost because you don't have to tune the entire set of model parameters. Full fine-tunes are more expensive because you are adapting parameters at every layer of the model's architecture. LoRA lets you change a smaller subset of weights and biases, for example 2-5% of the parameters, and get ~80% of the results. For testing purposes, this is fantastic, but for an end product, we found the results were not so great.

Moving to the larger Llama 2, 13B and repeating the process improved the results, however, it's not enough, which is why post training alignment is necessary (Stage 4 below).

## STAGE 3: INITIAL EVALUATION

Once this first training stage is complete, you need to work out what you've accomplished. Of course, you can just "plug the model in and ask it some questions" to quickly get a sense for the result - but a more comprehensive way to do this is using an evaluation or benchmarking tool.

This involves two parts. One is a set of benchmark questions or tasks that the model is asked or fed. The second is some way of evaluating the responses or results and scoring them. Most evaluations are fully programmatic these days. In other words, there is some sort of scoring model used to evaluate the results. This is fine, but once again, human evaluations are superior. While we've not built a product here, a Lightning-enabled evacuation tool has been on our roadmap for some time. This would allow the same community to rank or score responses (in turn with some sort of other, novel consensus mechanism) to earn Bitcoin and be paid out in real time for their contribution.

We have used a blend of internal human evaluation (ie; our team) along with some programmatic approaches.

I should also note that we have built the most comprehensive benchmark set of questions for a Bitcoin model to be tested against. It is 500 of the most important, common, pertinent and nuanced questions in Bitcoin. It is this set that we're training our model to perform well against and invite anybody else working on a Bitcoin-related model to come and test against this question set!

## STAGE 4: REINFORCEMENT LEARNING

Once this initial evaluation is complete, you get a sense for where the model is weak, where it is strong and where it needs further alignment. Reinforcement learning is a bit like fine-tuning but more dynamic. It's often called RLHF (Reinforcement Learning by Human Feedback) which has proven to be the most effective (once again, human involvement - sensing a theme?).

There is also RLAIF, which is similar but using purely model responses as the reinforcement examples. The latter is not as effective, but can be done faster. The key is to find a balance between both, because the RLHF process can quickly become expensive.

There are also different means of implementing the RL stage. Two of the more effective are PPO (Proximal Policy Optimisation) and more recently, DPO (Direct Preference Optimisation). With PPO, the learning is step-by-step, based on immediate feedback. In DPO, learning comes from comparing finished products and picking the better one. Both are effective and once again, the process is experimental.

Once again, this step of model training is ideal for a Lightning-enabled platform to engage human participants. We have actually built this and will roll the feature out shortly in our training app. It will primarily consist of the user ranking multiple answers to a single question. The best will then be used (in aggregate) to train a reinforcement model that will perform either PPO or DPO.

## STAGE 5: FINAL EVALUATION

Once some level of "alignment" is achieved, we reach the final evaluation stage. This is not so different from the preliminary eval, except that you are hoping to see a change in the results and a higher overall score.

If you've done a good job with the data, in the first place, and you add a little bit of luck into the magic that is "training" the model, you should come away with a positive result. This is not always the case, because small changes or mishaps upstream can turn into bigger issues later. But that's just the nature of the process - and why I call it more art than science at this stage. I am sure this will change as the industry matures, but I hope it's been made clear. We are in the very early days and much more experimentation needs to be done before it becomes more science than art.

In the next section we will review some of the Data, gathered up from the model training. There are some very interesting numbers here.

# 2. MODEL TRAINING DATA

In order to begin bridging the gap between the pain points outlined above, and the solutions a Bitcoin-native language model could provide. When going through the above process, we gathered, curated, generated, formatted, transformed, cleaned and validated a lot of data.

What follows are some statistics.

## 2.1 NAKAMOTO REPOSITORY

Anyone from around the world can contribute to this repository, which we envision evolving into the largest collective knowledge base of Bitcoin and Bitcoin-related content.

As it stands, anyone can search for anything and access the URL, or if they're available, the text file or PDF. A person browsing can use the various filters to search on a granular level to find what they are looking for. Individuals can also filter by content type (ie; book, directory, essay, article, blog, YouTube video or podcast), or format (link, txt, PDF, ePub) and search by author or the name of the specific piece of content.

We have a total of 33,533 contributions (see Fig. 8) in the repository currently, which is not only Bitcoin specific, but also contains roughly 14,000 items from the [Mises Institute](#) who graciously donated their entire database to the project.

It's worth noting again, that the quantum of data here is not as relevant as what we "do" with it. Quality is far more important than quantity. While the sheer number of resources is substantial, transforming it all into a format that is useful for the model is 90% or more of the challenge.

Our hope is that we can successfully do this, and embed the best of these ideas into the core model, then perhaps other people won't spend years having to filter for signal amongst the sea of noise out there.

## 2.2 HUMAN FEEDBACK

As described earlier, the human feedback component is the last mile, the 20%, that makes 80% of the difference by providing higher quality inputs for the model. Through the data pipeline I outlined earlier, we produced 53,000 question-and-answer pairs to be fed into the community tool.

The incentive for human participation, beyond the more altruistic benefit of contributing to a "Bitcoin project" such as this, is of course earning bitcoin as a reward. We've had approximately 300 people from different continents, including Africa, Latin America, Europe, Asia, Australia and North America, all participate.

In total, almost 1M Satoshis of micropayments have been auto-paid in this process. The number is not high because at the time of this report, the rollout of actual payouts was fresh. Through our broader bounty program, manual payouts for developer assistance and data transformation has exceeded 40M Satoshis.

For the purposes of this report, we will look into the automated component as this holds the greatest potential. As noted earlier, there are three core data-related features inside our tools available to participants. Let's look at the data from each now.

## DON'T TRUST, VERIFY (DTV)

This feature presents users with a question/response pair, along with several possible answers they can respond with. They can keep the question and response, discard the question and response, or edit the question and/or response. As of December 1st, 2023, there have been a total 43,663 unique responses.

The consensus mechanism requires a threshold of users to perform the same 'final' action (ie; keep or discard). On successful consensus, the data is moved to the next stage of the pipeline and removed from the feed. Simultaneously, points (Sats) are available in the accounts of the participants who were in consensus.

The edit function creates a new data artifact (updated Q&A) and places that at the top of the DTV funnel, for others to now "keep" or "discard", while the original remains in the funnel for others to keep, discard or edit. This continues until good data is ultimately moved on and bad data is cleaned out.

Does it create multiple variations of the same or similar questions and answers? Yes. And for training purposes specifically, this is a good thing. We designed it so the process kills two birds with one stone.

Some stats below:

→ Total responses: **16,180**
→ Answers accepted: **7,114**
→ Answers discarded: **619**
→ Answers edited: **2,317**
→ Pending consensus: **8.436**



**NAKAMOTO REPOSITORY ITEMS**

- Course (1)
- Podcast (4)
- Other (10)
- Directory (16)
- YouTube Playlist (61)
- Essay/Article (236)
- Book (433)
- YouTube Video (692)
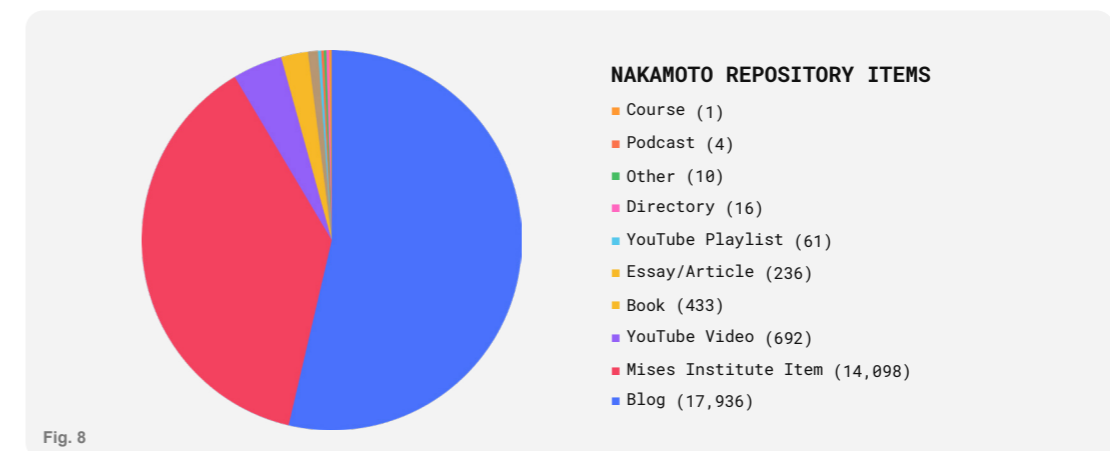- Mises Institute Item (14,098)
- Blog (17,936)

Fig. 8

## WE ARE ALL SATOSHI (WAAS)

For We Are All Satoshi, users are presented with a question and a blank slate to respond to. The idea is to "pretend" you are Satoshi (not the real person, but this AI model), and answer it in the way you would imagine such a question to be answered.

This requires individuals to come up with original answers, which obviously takes more time, and as such, rewards 10x the number of Sats than the DTV feature does. Similar to the "edit" function in DTV, new responses are placed in the top of the DTV funnel for people to verify, via accepting, discarding or editing.

This feature has provided valuable, novel and unique human-generated data for use in the model. Some statistics:

### Stats here

→ Total responses: **1,773**
→ Answers accepted: **1,013**
→ Answers discarded: **40**
→ Answers edited: **512**
→ Pending consensus: **213**

### FUD BUSTER (FUD)

The FUD Buster feature is basically a clone of WAAS, but instead of just being presented with another question, users are presented with either a statement or specific question implying some sort of FUD (fear, uncertainty and doubt) regarding Bitcoin.

Once again, they can respond with up to 2,100 characters. These replies are fed into DTV until consensus is reached, ie; the final data artifact is either accepted or rejected.

As you can see below, this has less participation from the cohort, and we imagine that is because it takes longer to think through an answer than it does to validate.

→ Total responses: **1,529**
→ Answers accepted: **1,017**
→ Answers discarded: **18**
→ Answers edited: **354**
→ Pending consensus: **140**

## OTHER FEATURES

It's worth noting briefly the other features we developed, which can at some point include, integrate or use Bitcoin / Lightning in some way.

### LOGIN

We have three methods for sign up and login, which are all essentially the same flow:

→ Nostr
→ LNURL, or
→ Email

Email with a magic link still seems to be the most popular, but the fastest is certainly LNURL because you just scan the QR code and you're in. The Nostr flow is similar to the email flow, in that you input your NPUB + a relay, and we send you a magic link / 6-digit

code that you can input and sign up / log in. By the time you're reading this, we should have also rolled out NIP-07 so that people can quickly log in with their browser extension (eg; Alby).

We encourage people to associate a few credentials to their profile, in case they lose access to one. In this way, there are multiple back ups for log in, which we believe is very useful.

We hope that at some stage, tools like Slashtags from the team at Synonym, will roll out and enable other ways for people to "log in" while owning their credentials. Likewise, in time, we will seek to integrate these logins more deeply with other Bitcoin, Nostr, etc. native features.

### PAYOUTS

This is simple and straightforward. For ease of setup, we partnered with OpenNode for LNURL withdrawals. Yes, the training app is currently custodial, but that's more a function of the bandwidth we have available for development right now. A non-custodial solution is definitely on the roadmap, so that payments flow directly to users. For the moment, since we are the entity doing payouts to contributors, it makes sense to have a custodial set up.

At the time of this writing, 58,344 points and 594,380 sats have been awarded for 43,663 contributions, from 296 contributors. These points translate to Satoshis. We used points in between, from the outset, so we could more easily boost conversion rates (from points to Sats) for example if we want to do a "bonus" for active users, or targets, or the leaderboard.

Individuals can see the breakdown of the points they've earned on their personal "stats" page, in their profile (see Fig. 9).

### STATISTICS

We built a simple, but very visual "stats" page, as mentioned earlier, so people can see how much they've earned, how many contributions they've made, the number which have been kept or discarded, and more.
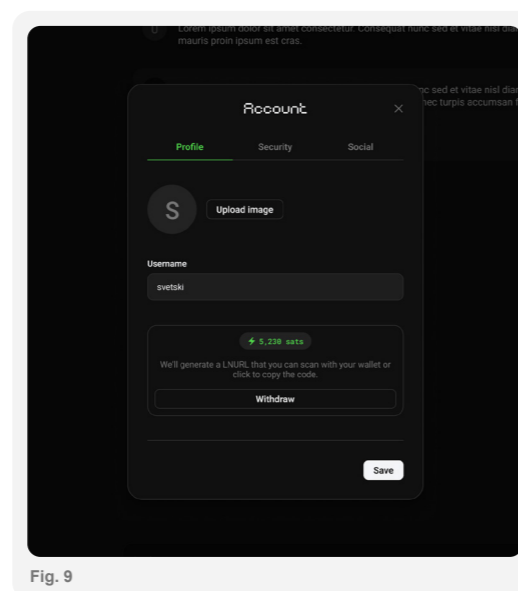
Fig. 9

As we roll out more features and mature the product further, we aim to build more features into the stats page and integrate it more deeply with our leaderboard and other social features.

### LEADERBOARD

Of course, even stronger of a driving factor than straight incentives and internal statistics, is the human need to compete. We built a simple leaderboard (see Fig. 10), whose current functionality is more focused on total "points" earned, in order to give contributors a sense of what others are doing and drive them to compete.

The leaderboard allows for gamification, bonuses and benchmarking amongst the community, and in time, we will look to integrate other features into it that are more social in nature. More on that below.

You can check out the current leaderboard here and sign up to participate: https://www.train.spiritofsatoshi.ai/app/leaderboard

### SOCIAL FEATURES

Finally, one of the most powerful benefits you realize when you combine a tool like this with an internet-native monetary protocol and decentralized identity/communications protocol like Nostr, is the ability to make the whole experience more social, interconnected and rewarding.

To begin with, one low-hanging fruit would be to associate contributors' accounts with their Nostr and Twitter accounts so you can follow these people, but in time, the more interesting features will include the ability to:

→ View contributions from other users.
→ Highlight those contributions and share across Nostr-related apps.
→ The ability to Zap other users for status, contributions and more.

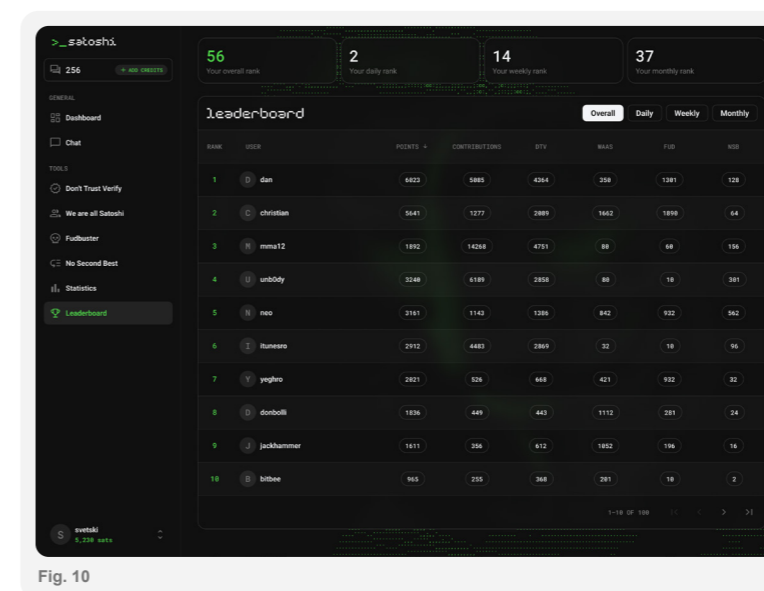As we build this tool out further, these will be areas of great interest.

Fig. 10

# 3. MODEL TRAINING FUNDAMENTALS

**A**s should be evident now, training a large language model is no small task. Researchers and companies have been testing out new techniques and developing an understanding of what yields better results, but it's still early days. Like I've said, it is very much both an art and a science. Perhaps more so in our case, since the fundamental paradigm the Satoshi models are meant to reflect are opposite or out of phase with just about every kind of model out there right now, whether open or closed source. How we dealt with these challenges was discussed earlier in the report.

In this section, I'd like to help elucidate the differences between training, tuning and augmentation - which is often mis-referenced as a way to "train your own model".

## TRAINING VERSUS FINE TUNING VERSUS AUGMENTATION

There is a lot of misinformation out there regarding "training" models. If you trust the average 'AI bro' on Twitter, you'll believe that you can just upload a few PDF's, a couple of books, a podcast or company's financial report, to "train your own AI". This is completely false.

For a number of reasons, two of them being:

→ You cannot train a model (or even effectively fine-tune it) without a sizeable corpus of data.
→ Even if you have A LOT of data, the data must be transformed into a format which represents the style and kind of output you want the model to produce. For example, if you're training a model to answer questions, you cannot just "feed it some books." You actually have to transform the content within the book into a set of questions and answers, which become the examples, or "training data set".

This was made abundantly clear in Section 2 of this report. In the following section, we will explain what this all means, and the different options one has with respect to training a model, fine-tuning it, or augmenting an existing model with a semantic database that it references.

### UNDERSTANDING LARGE LANGUAGE MODEL (LLM) TRAINING

Large Language Models are AI models that can understand and generate human-like text. Unlike more focused AI applications, LLM training doesn't target a specific domain or task; instead, it aims to develop a comprehensive understanding of language, often across multiple languages and contexts. This training involves feeding vast amounts of quality* text data into a model, enabling it to learn patterns, nuances, and structures of language.

Notice the use of the word "quality" here. Quality is a subjective term. In the context of data for LLM training, quality refers to how representative the training data set is, with what you want the model to output later.

The process of LLM training is complex and multifaceted. It's not just about accumulating data but also about preparing it in a way that's conducive to learning the intricacies of language.

Imagine it as teaching a child language by exposing them to an extensive library of books, conversations, and writings, but at a scale and speed that's only possible in the digital world.

### THE ESSENCE OF LLM TRAINING: COMPREHENSIVE LANGUAGE UNDERSTANDING

LLM training is about creating a foundation of language understanding that's broad and deep.

The trained model should be capable of not just repeating what it has seen but also generating new, coherent, and contextually-appropriate content.

The Process and Components:

→ **Data Collection and Curation:** The first step involves gathering a diverse and extensive dataset. This dataset can range from general literature, websites, and articles to more specific texts depending on the desired expertise of the model.
→ **Data Preparation:** The collected data needs to be cleaned and formatted. This involves removing irrelevant or sensitive information, correcting errors, and ensuring that the data is in a uniform format for the model to not only process efficiently, but as mentioned earlier, representative of how you'd like the model to output content later.
→ **Training Phase:** In this phase, the model is exposed to the prepared data. Using various algorithms, the model learns to predict the next word in a sentence, understand context, and generate coherent responses, all through developing relationships between letters, words and sentences. The weighting and biases between these elements are the "parameters". This phase requires large amounts of computational power and can take days, weeks or even months depending on the model's complexity and the size of the dataset.
→ **Dataset Size and Variety:** The effectiveness of an LLM is directly related to the size, the diversity or specificity of its training dataset. A larger and more varied dataset enables the model to develop a more nuanced understanding of language and its many applications, but also leads it to produce more generalizations. In other words, the more general you want the capabilities of the model to be, the more varied your dataset must be (this is known as your data-blend). The more specific the capabilities, the more specific the dataset needs to be.

### COMPUTE REQUIREMENTS

The computational requirements for training an LLM are not trivial. High-end GPUs or TPUs, often available only to well-funded organizations or research institutions are needed. And of course, this all requires energy. The process is not only data-intensive, but compute-intensive. At the micro level, this is not a concern. Companies like us just plug into what's available, use "free credits" from cloud providers where possible, and seek compute through whatever means possible (centralized or decentralized like GPUtopia).

Of course, at the macro level, this is a concern, and it will be interesting to see what the Bitcoin industry can teach the AI industry when it comes to the efficient scaling of compute.

All this is to say that when people tell you that "you can train a model on your own data" - they have absolutely no idea what they're talking about.

### UNDERSTANDING FINE-TUNING IN AI MODELS

Fine-tuning particularly follows the pre-training phase of the LLM development cycle, although in computational terms, it's quite similar. The difference here is specificity of the data and the time involved. It's akin to honing a broadly-educated mind to specialize in a specific field. After an LLM has been pre-trained on a vast, general dataset to understand language, fine-tuning adjusts the model to excel in specific tasks or comprehend particular domains.

Think of fine-tuning as customizing an all-purpose tool to perform specific jobs with greater efficiency and accuracy.

This phase is crucial for tailoring a model to specific needs, whether it's understanding medical terminology, generating marketing content, engaging in casual conversation, or in our case, speaking like a bitcoiner!

### THE ESSENCE OF FINE-TUNING: SPECIALIZATION OVER GENERALIZATION

Fine-tuning shifts the focus from a general understanding of language to specialized knowledge or capabilities. It actually involves retraining the model, but now with a dataset that's closely aligned with the intended application.

### THE PROCESS AND VARIATIONS

There are three "general" categories for fine tuning. Once again, this is not 100% the case, all of the time. It's just a useful way to understand it.

→ **Full Fine-Tune:** This involves retraining the entire model on a new, domain-specific dataset. It's like giving the model an intensive course in a new subject, reshaping its understanding and response patterns to align with specific requirements.

→ **Low Rank Adaptation (LoRA):** LoRA is a more targeted approach to fine-tuning. Instead of retraining the whole model, LoRA adjusts only a small fraction of the model's parameters (usually the top layers). This method is efficient and requires less computational power. It's particularly useful for fine-tuning models where access to the entire model structure is restricted or when computational resources are limited. It's a bit like an 80/20 rule, but more like 80/2. You tune the 2% that matters, to give you 80% of the result.
→ **Partial Fine-Tuning:** In some cases, the fine-tuning process might be constrained to the top layers of the model. This form of fine-tuning still allows significant customization but within the framework of the original model's broader understanding. Fine-tuning OpenAI's models is such an example.

### APPLICATIONS AND VARIED APPROACHES

The choice between full fine-tuning, LORA, and partial fine-tuning depends on several factors:

→ **Intended Use:** The specific task or domain for which the model is being fine-tuned can dictate the depth and approach of fine-tuning.
→ **Resource Availability:** Full fine-tuning requires substantial computational resources and data, whereas LORA is more resource-efficient.
→ **Model Accessibility:** Some models, especially proprietary ones like DaVinci by OpenAI, may have limitations on how deeply they can be fine-tuned.
→ **Testing:** If you're testing, it's often best to start with a LoRA tune and then if the results are positive, move onto a full fine-tune.

### UNDERSTANDING REINFORCEMENT LEARNING

The final stage in LLM development focuses on aligning the model with specific human standards and preferences. This 'last mile' stage is essential for refining the model's decision-making capabilities and ensuring its outputs align with desired outcomes, particularly in terms of relevance, style and accuracy.

Imagine this as the final tuning of a high-performance engine, ensuring it not only runs smoothly but also responds precisely as intended. There are two main options for reinforcement learning, and a blend of both can be used. Let's look at each.

### REINFORCEMENT LEARNING FROM HUMAN FEEDBACK (RLHF)

RLHF requires human feedback to build a reward model, in order to then further LLM refinement. The steps are:

→ **Collecting Human Comparisons:** RLHF starts with human evaluators providing qualitative feedback on the model's outputs, effectively teaching the model what is considered a desirable response. Think "ranking" and "scoring" response variants.

→ **Training a Reward Model:** This feedback is used to train a separate 'reward model'. This model learns to predict which responses will be favored based on human evaluations.

→ **Iterative Refinement:** The LLM is then fine-tuned using reinforcement learning techniques, such as Proximal Policy Optimization (PPO), or the more recent DPO, to maximize the rewards as predicted by the reward model. This process is iterative, continually evolving the model's output quality based on new feedback.

### REINFORCEMENT LEARNING FROM AI FEEDBACK (RLAIF)

Very similar to RLHF, except we use existing LLMs to get the feedback. This is of ultimately lower quality, but also lower cost in time and money.

→ **Automating Feedback:** RLAIF involves using another AI model to provide feedback, making the process more scalable. This AI-generated feedback aims to mimic human evaluations, guiding the LLM towards desirable outputs.

→ **Challenges in RLAIF:** While RLAIF enhances scalability, it also introduces complexities in ensuring the AI feedback's quality and reliability, which is crucial for the model's accurate and ethical alignment.

### CONCLUSION

Reinforcement Learning, whether RLHF, RLAIF or some blend, plays a vital role in the final stages of LLM development. It is the last mile alignment stage, and really puts the icing on the cake, so to speak.

Now that we have the training stages out of the way, let's look at what most people erroneously call "training" today, and understand why it is fundamentally different to training a model, but still useful in particular contexts.

## AUGMENTATION

Retrieval Augmented Generation (RAG) is the most popular way to augment or enhance a model, so we'll put our focus here. RAG is a novel way to get a model to produce responses that are more accurate or "relevant" to a domain or point of view. Contrary to popular belief, RAG has nothing to do with actually "training a model". Instead, it focuses on augmenting the capabilities of an existing model. This augmentation is achieved not through retraining with new data, but by enhancing its responses through a sophisticated use of embeddings and external data retrieval.

Think of it as a smart way to do dynamic prompting by abstracting away the context injections using semantic tooling.

It's a bit like asking a model to answer a question by referencing some specific context you pasted into the prompt. Imagine you just copied a relevant section from a book, pasted it into ChatGPT, then asked the model to answer a question by referencing that context. It's actually pretty simple, conceptually speaking.

In fact, most people who use ChatGPT (or any other model) do this already, only somewhat manually. They make sophisticated prompts so that the model can reply more accurately. RAG just allows you to do it dynamically and programmatically. It abstracts away the manual process.

### THE ESSENCE OF RAG: AUGMENTATION OVER TRAINING

RAG enhances AI applications by allowing them to dynamically access and incorporate information from external databases. This method effectively broadens the AI's knowledge base without altering its foundational training. The core AI model, already trained on a substantial dataset, is coupled with a retrieval system that fetches relevant information from a vast external database in response to specific queries or content requirements.

### THE PROCESS AND COMPONENTS

→ **Data Ingestion and Embedding:** RAG starts with embedding large amounts of data into a vector database. This database is optimized for quick semantic searches, crucial for retrieving relevant information rapidly. Embedding is simply the process of transforming data into a vector form that can be efficiently processed, read and referenced by an LLM.

→ **Contextual Response Generation:** When a user query is received, RAG identifies relevant data from the vector database and uses this context to enhance the AI model's response. This process involves interpreting the user's input, searching for pertinent information, and then integrating that information into the response.

→ **Size of Dataset & Vector Space:** The beauty of RAG is that you can augment a model on a small dataset, for example a single book or article, or you can use massive datasets, although that requires a lot more upfront work with data chunking, metadata and, assuming you want high quality results, ensuring all of the embeddings are of high quality (this can take quite a bit of time).

### PRACTICAL APPLICATIONS ND LIMITATIONS

RAG is valuable in areas where AI responses need to be supplemented with up-to-date or specialized information. However, RAG's effectiveness hinges on the quality of the external data sources and the system's ability to accurately match query embeddings with relevant information.

The complexity of setting up and maintaining such a system, especially in dealing with vast and continually updating data sources is precisely where things get challenging.

Furthermore, the core model has not been changed, so it's not producing anything novel or unique. The model is still the same underlying model, and as soon as a question is asked that's outside of what's in the vector store, it will revert to default, or not answer.

If your goal is a little widget, this is a useful solution. If your goal is to reference an internal document more easily, or perhaps turn your company FAQ's into something that you can reference conversationally, then great. But this is not a new model, and it will not perform as well as a fully trained model will.

### CONCLUSION

In summary, LLM training is a powerful but resource-intensive process aimed at creating AI models with a broad and deep understanding of human language.

It's a complex endeavor that combines data science, machine learning, and linguistic expertise, resulting in models that can interpret and generate human-like text across various contexts and applications. The training process not only shapes the capabilities of the model but also sets the limitations within which it operates.

Fine-tuning allows for the customization of a general model to meet specific needs and perform specialized tasks. It can be done via a low-resource approach like LoRA or as a full update to the model's parameters. Reinforcement learning is the last mile of the process, and aligns the model.

Finally RAG, or other approaches to augmentation, are not training, but enhancements or wrappers on models which are great for very narrow applications, prototyping and demonstrations.

Understanding these different elements and their implications is key to leveraging the full potential of AI in a targeted and efficient manner.

# 4. CROWD-SOURCING, SPECIALIZED, OPEN SOURCE, LLMS

In the 60s, businesses adopted the phone. In the 2000's, they built a website. In the 2020's and beyond, they will all have AI agents. If the website is the "store front", the AI avatar will become the "digital employee". Always on, always available, always reliable, never complaining.

## HYPOTHESIS

As this technology develops, as the price of compute comes down, the demand for higher fidelity avatars and agents is likely to increase.

More companies, industries and brands will look to have their own bespoke agents represent them, and perhaps on a long enough time scale, individual agents for individual people. Large, general models won't cut it for such applications. They are good for general consumer use-cases, but they cannot match the performance, accuracy and price of specialized models in domain specific tasks.

We are already seeing evidence of this with the number of models on Hugging Face and the countless other enterprise grade, private models that companies like Mosaic are building for large clients today. In the next 10 years, it's likely that every company, every brand, every country, city, influencer, CEO and small ``Gig entrepreneur'' will have their own AI assistant or avatar. And I do not mean some RAG model using Open AI, nor even a LoRA fine tune of an open source model. I mean high-fidelity LLM driven agents tied to 21st century knowledge bases, that can accurately represent someone, or some brand or some point of view.

Building such high-fidelity models is possible today (mostly), but it's extremely expensive, and reserved for the Google, OpenAI and Meta's of the world. See the new avatars that Meta spent hundreds of millions building for the Messenger Chat. To make this technology more widely accessible, three things need to happen:

→ Compute must come down in price, significantly.
→ The data available must be transformed into a useful format.
→ The frameworks, pipelines and tools to do the above need to be built.

It is the last two where crowd-sourcing model development is likely going to have the biggest impact, particularly for domain specific models - and perhaps one day, even larger generalized models.

## INCENTIVISING THE CROWD

If the hypothesis is accurate, then what we outlined earlier in the report suggests that human feedback ) see Fig. ??), data and involvement is going to need to scale up tremendously to meet this demand. This is extremely difficult to do with fiat money rails and the legacy banking and payments networks, especially at a reasonable cost.

This is where Bitcoin and micropayment networks like Lightning come into play, along with user-owned accounts or identity, as with Nostr. They are global, accessible to anyone with an internet connection, platform and application agnostic, and contrary to what some might say, are very easy to use.

There needs to be a way to connect the people who have the time and knowledge, but not the financial resources, or access, with the groups, companies, communities or projects that have funding, but neither the time to generate or curate data, or in many cases, lack the budget of an OpenAI, and must therefore seek cheaper labor.

Imagine an American company or community trying to do micropayments for data curation, creation and verification, for workers all over the world, in Latin America, Africa, South East Asia, the Middle East. It's impossible. Upwork doesn't cut it. These locations are full of latent talent, all who have no access to banking or payment rails. Couple the fact that "data is the oil of the 21st century" and as has been seen, is the determinant of model quality, then those who can find a way to leverage this talent pool and effectively incentivise them, will have a major advantage.

It's for this reason we're expanding upon our lightning enabled crowd-sourcing tool. We built this for ourselves, to solve our own problem, ie; Building a Bitcoin model. Turns out the Bitcoin model suite, while useful, is more of a public utility with a commercial application with just the Bitcoin space. On the other hand, the tool that helped us leverage the community to build it, is perhaps one of the most important tools for the broader AI and data industry.

## DOMAIN SPECIFICITY

A final note on this thesis. More important than just having open source models, is the existence of multiple models for different use-cases and domains. Whether closed or open source, more flavors and more variants equals more actual decentralization and a more free market.

If the future is going to be "multi-model" and if "everybody will have their own AI/agent" then the frameworks for making this possible must be built. Specialized AI models, be they medical, legal, recruitment, finance and investing, learning and development, training, education, specialized customer support, analytics, operational support agents or just characters - will need the power of the crowd.

If you're a company, content creator, brand, industry body or community interested in developing your own bespoke model, and have a community, crowd, or customer base that you can leverage, please reach out to us. Our platform enables the development and deployment of specialized AI models for any industry vertical. We can help you:

→ Unlock the real value of your data. Transform raw data into a series of formats, useful for either model training, or semantic storage so that it is language-model-readable.
→ Train your model efficiently and effectively. Training a model is currently more art than science. Our process ensures that you spend less time experimenting and more time building.
→ Human Feedback and Reinforcement Learning. Automated pipelines are critical, but the last mile needs humans. Our lightning-enabled tools enable community members, employees, researchers or anyone with domain specific insight to help, irrespective of their location, geography or banking set up.
→ Deploy your model effectively and efficiently. There are hundreds of deployment solutions available. We help you analyze and determine the best fit for training, deployment, and ongoing inference.
→ Connect your model to other tools / data sources. The real power of AI will one day lie in its ability to use tools and multiple sources of data. Our framework makes this possible through "tool-training" and element connectors.
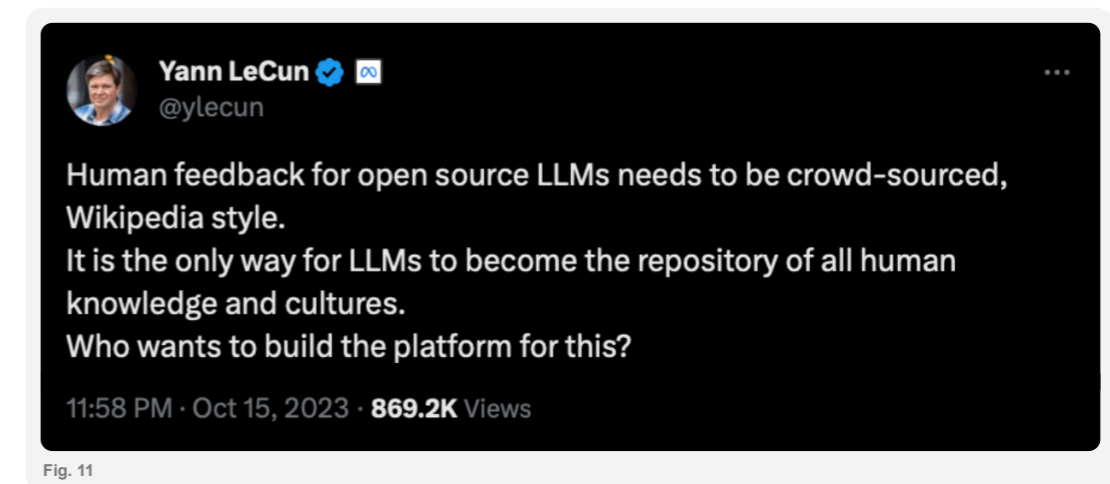


> **Yann LeCun** ✓ ∞
> @ylecun
>
> Human feedback for open source LLMs needs to be crowd-sourced, Wikipedia style.
> It is the only way for LLMs to become the repository of all human knowledge and cultures.
> Who wants to build the platform for this?
>
> 11:58 PM · Oct 15, 2023 · **869.2K** Views

Fig. 11

# 5. MYTHS IN AI

**T**here are myths and misconceptions in every industry. The AI space is near the top, due to its potential, and the poorly understood nature of intelligence. People are either over or underestimating things, and simply misconstruing capabilities, every single day. This section is a rapid fire myth busting section for easy future reference.

## YOU JUST NEED MOAR DATA

You often hear that ChatGPT was trained on the entire internet. This is not only false, but off by so many orders of magnitude, it would make your head spin. GPT-3 -175B was ~600 GB (300B tokens) while the internet was 64,000,000,000,000 GB. That's a bit like comparing your neighborhood to the size of the Sun. Or more accurately, if all the data on the internet was represented by the entire surface of the Earth, then all of ChatGPT's data would only be represented by about 478 square centimeters (or about 74 square inches), or approximately the area taken up by a typical dinner plate.

**Why is that so?**

It's because most of the data out there is not in a useful format for training a language model. In fact, you can think of data like untapped, raw materials: it has to be cleaned and refined, before it can be used.

Remember that Language Models are trained on the relationship between words and sentences. The examples must be representative of what you want the model to produce. Less is more. Higher quality, and only if possible, greater quantity - but not at the expense of quality.

## YOU CAN JUST "TRAIN A MODEL ON YOUR OWN DATA"

Well, not really. When people are saying this, what they're often talking about is some form of retrieval augmentation, ie: "RAG".

This is not the same as fine-tuning or training a large language model! RAG stands for Retrieval Augmented Generation, which utilizes external data sources at inference time to add additional context for responses. This is a very different process and technique than pre-training or fine-tuning a model. RAG should be treated as a supplemental tool, but it is not actually "training". Nothing is changed about the underlying model at all. It is a way to prototype, or create a "chat to your documents" mini-agent, but it's not "your own model" and it's definitely not robust enough to have long, meaningful or contextual conversations with.

RAG shortens the context window because of all the extra context you need to inject into the abstracted prompt. This means that after a few responses it will lose context. You can mitigate this with a sliding context window prompt, but this is not a great fix assuming you want a fluid, useful dialogue with the model.

Ultimately, training your own model on your own data will require a lot more data than what you would use in a RAG scenario, and is an entirely different exercise which remains expensive and out of reach for the average person.

## WE NEED "UNBIASED" LANGUAGE MODELS

Bias is not something that can or even should be removed from language, discourse or personalities. Bias is another word for preference, or opinion, or "worldview." All discourse, all data, all information has within it an implicit bias.

When it comes to language models, since they reflect some aggregate of the data they're trained on, they will fundamentally also reflect that bias. It's inescapable. The workaround is of course to put guardrails on the model, to inject pre and post-framing for every response (as is done with ChatGPT these days) but that doesn't remove the bias - it just creates a bad user experience.

Trying to eliminate bias is like trying to flatten everything. It's a Quixotic pursuit. The focus instead should be on being clear about what the bias is, and building many alternatives. Since a bias is just a model of the world, we want many of these, not just one or a few.

When people are talking about "unbiased AI" they are either misinformed, naive, or in some cases, using that as a way to claim a moral high ground in order to impose rules around what language or styling is "acceptable."

## AI GETS RID OF THE NEED FOR HUMAN WORK

The idea that an AI will one day replace humans — either by taking your job or by annihilating humanity — is a scary concept. It has inspired a plethora of films and books, so nervousness about the implications of AI is understandable.

But if the last few years have demonstrated anything, it's that when people are scared and falsely-informed, that they make the worst decisions.

It's important to note, models not only perform significantly better when there is human feedback and human-generated data involved, as per what we've documented in this report, but that models are only as useful as the person who uses them. Nothing has changed about the nature of tools. There is an actor and a tool. AI is merely a tool which if used well, can yield superior results.

## AGI IS AROUND THE CORNER

Related to the above unfounded fear is AGI. It's nebulous enough to be scary, and people, who otherwise have not enquired into the nature of either intelligence or consciousness, often believe that somehow something sentient will emerge from the circuits. As a result, we must either ban, or form a regulatory body to "manage it", of course, "for our safety".

I am personally not of the opinion (which is not shared by everyone) that AI is suddenly going to become sentient and rule or take over the world. AGI and the singularity is a red-herring.

The real danger is of AI as a tool being wielded only by those who have questionable intentions, or a poor track record with other tools. Examples abound.

The existential risk is that such entities or groups embed power AI tools into every layer of society, and therefore reduce human liberties and dull the color of life.

It is this threat we want to counteract. The idea is to build AI-enabled tools that enhance human flourishing, and bring more color and nuance to life.

## "WE'RE ALL GOING TO HAVE OUR OWN AI, ON OUR LOCAL MACHINES"

This may happen, one day. But not for a while. Perhaps even decades. Why? First of all, it's related to what's mentioned above about fidelity. The technology has a long way to go before it becomes more science than art. To properly train and tune smaller scale models for every person will require a whole host of frameworks, pipelines and tools that simply do not exist today. Furthermore, the compute necessary is just not available.

Second, and this is a less understood, more insurmountable factor - as better large-scale cloud models are released, they will raise the minimum acceptable bar, and thus make these smaller DIY models less interesting and useful. This has more to do with the human condition than it does the efficacy of the self-hosted models.

Notice how the first time something happens, "it's a miracle" and then it becomes normalized. It happens with flying on an airplane, with using the internet and it happened with ChatGPT. Everyone lost their minds for a minute, and now it's just another app.

How this relates to the point is that at no point, no matter how much better compute and local hosted models get, will they exceed the capabilities of larger, cloud-hosted models. And these larger cloud-hosted models (whether ChatGPT or other) will set the bar for usage, quality, functionality, etc. Using your local model will be like going from an airline back to a wooden sailboat, or from a car back to the horse and buggy.

Now, before you say: "but there are small models outperforming the large one's already", please read the next point.

## BENCHMARKING AND EVALUATIONS

This one is not so much a myth, but a misconception When people see that "x" model has outcompeted "y" model with more or less parameters, they immediately assume x model is better. This is not necessarily true.

Why? Because, evaluations and benchmarks are not only subjective, but they are narrow and can only evaluate models within the window they apply. So what happens is two things:

→ People game the results to hit the leaderboards. Models are often tuned specifically around a series of evaluation metrics and benchmarks. This means they perform well in those tests but not so well outside it.

→ This creates the false assumption that the models are broadly better than they really are. It's easy to go look at a Hugging Face leaderboard and assume that it applies across everything. This is why GPT-4 continues to outperform all of the open source models, and why everyone continues to use it.

This is not to say benchmarking and evaluation is bad. It's just misunderstood, and as a result, people project forward erroneously. In fact, benchmarking and evaluation is necessary, particularly for projects like ours, which are domain-specific. Because we can, within our domain, show that what we've built, outcompetes models x, y and z.

We cannot claim our model is useful outside of this context - but that is fine, because we're not claiming anything beyond that.

## OPEN-SOURCE VS CLOSED-SOURCE

Also not a myth, but a series of misconceptions.

The first confusion relates to what is being open-sourced? The data or the model? Notice that very, very few groups open-source the data sets. In fact, I'm not sure I've seen one major "Open-Source" model, also open source their full database. This is because it comes with a whole host of legal implications.

What they do willingly open source are the weights and biases. And this is great, but outside of a few data scientists around the world, it doesn't mean a lot to most people. Very few are going to print out the parameters and check for themselves.

This is not to disparage open source at all. It's extremely important because so long as some people can check, that is great. Therefore in this context, closed source is not so different. It basically means you don't know about what you wouldn't understand anyway. In other words, there is nothing you would do with the weights and biases anyway.

But…and very importantly - where Open Source shines is that it enables anybody, anywhere to take the current model (depending on the OS License) and adapt it. They can re-train, fine-tune and really turn it into something new. That's precisely what we've done with the Satoshi suite of models, and the upcoming "Code-Satoshi".

The most important thing once again is application and honesty about "does it do what it says on the label"? In other words, if you want to build something more proprietary, just tell people what it does and do not pretend it is "unbiased." This is once again where Open Source does shine, because there are a few great magicians out there who can check whether the ingredients are really there, and can bring such things to light.

The final note here is on crowd-sourcing. If we can successfully work out how to build these models with the help of the crowd, then of course they should be fully open sourced, and become "utilities", so to speak. This is our mission with the Satoshi suite of models, and Max Webster from VC firm Hivemind ventures discussed in his essay earlier this year. He specifically wrote about ways that Bitcoin and the Lightning Network can power open source models to win.

As bitcoiners, we appreciate the open source nature of Bitcoin and the Lightning Network code. It's a fantastic read, as is the following post from the team at Turing Post: https://www.turingpost.com/p/openvsclosed

## CONCLUSION

**W**e hope this report has provided insight, education, and value to you. If it has, please feel free to share it around and together we can help people understand this technology better.

### WHAT'S NEXT?

We aim to roll out a full suite of Satoshi Models, of different sizes, for public access on our website, and for download from Hugging Face. At the time of this writing, we are still in the development stage, but perhaps by the time you read this, some or all of what we're planning will be available.

We will continue testing the models, tuning new versions, releasing upgrades and the like. We also plan to engage further with Bitcoin businesses to collaborate on product opportunities that can serve their customers.

We also plan to release at minimum, one of these reports annually - but considering the pace of change, perhaps we will need two! At the very least we will follow this report up with a second that includes a deeper exploration of the solution space and partnership opportunities with Bitcoin companies interested in deploying AI tools.

It's still too early to say how much AI, and in particular language models, will change the world. Whether it will have the size of impact that some say remains to be seen - but I am pretty confident that once the hype dies down, we will, over the coming decade, find clear applications and uses for such a tool.

In the meantime, it's important to step outside of the hype and critically-analyze what is and is not useful. It's very easy to get caught up in "potential" applications, and allow the imagination to go all exponential on you. It's a very human thing. Turning that imagination into something tangible is what a business and an entrepreneur does. It's our hope this report will be useful along that path.

Thank you for taking the time to read this, and if you're interested in collaborating or finding out more about ways in which AI tools can enhance your business, please see the following pages for further information and a sneak peak of what we're working on.

## ALEKSANDAR SVETSKI.
**& The Spirit of Satoshi Team**

# FOLLOW THE JOURNEY

## PARTNERING WITH SPIRIT OF SATOSHI

**A**s the next billion people worldwide ask the question, "What is Bitcoin?", it is essential that the answer is accurate, clear and engaging. We envision Spirit of Satoshi being both a character people can turn to to ask those questions, but also a tool for Bitcoin companies to better deliver this message and service customers at scale.

Several ideas came from the product discovery sessions we undertook, and we're actively honing in on which of those make the most sense to put resources toward.

If you are interested in being involved and building out a specific solution, whether as a Bitcoin business or a content creator, please reach out.

Likewise if you'd like to get access to an API to use Satoshi in your stack or your app. For inspiration, below is a list of products we're working on.

→ **Satoshi Language Model.**
  › Bitcoin Customer Success Agent. Using up-to-date information about your company, Satoshi can assist new users with onboarding onto your product, or learning about Bitcoin.
  › Content-Generation Assistant. Need help generating Bitcoin content? Don't we all. Use Satoshi for Twitter, newsletters, ideas for blogs and even scripts for new content.
  › Bitcoin Tutor / Guide. There is a new wave of Bitcoin-education companies on the rise. Satoshi can float on screen for students going through these courses, and even act as an "assessment agent" to help make the assessment more than just a "multiple choice" exercise.

→ **Code-Satoshi.** Our most exciting suite of Satoshi models are Bitcoin-coding assistants. The first version will specialize in Miniscript, but we intend to add much more in the coming months, including support for Liquid, RGB, BitVM and more.

→ **Crowd-Sourced AI Model.** If you're interested in building your own model, and would like to leverage what we built with the Lightning-enabled training tool, we'd love to help.

This is applicable well beyond the Bitcoin space, so if you've been thinking about building something in a domain that is not so mainstream, or perhaps just specific (eg: Self Defense, Homeschooling, Recruitment, Bitcoin Education, The Bible), please reach out to work with us. As outlined in Section 6 of the report, we can help you:

→ Unlock the Real Value of Your Data. Transform raw data into a series of formats, useful for either model training, or semantic storage so that it is language-model-readable.

→ Train Your Model Efficiently and Effectively. Training a model is currently more art than science. Our process ensures that you spend less time experimenting and more time building.

→ Human Feedback and Reinforcement Learning. Automated pipelines are critical, but the last mile needs humans. Our lightning-enabled tools enable community members, employees, researchers or anyone with domain specific insight to help, irrespective of their location, geography or banking set up.

→ Deploy Your Model Effectively and Efficiently. There are hundreds of deployment solutions available. We help you analyze and determine the best fit for training, deployment, and ongoing inference.

→ Connect Your Model to Other Tools / Data Sources. The real power of AI will one day lie in its ability to use tools and multiple sources of data. Our framework makes this possible through "tool-training" and element connectors.

Make sure you follow Spirit of Satoshi on Twitter, Nostr and LinkedIn to stay up-to-date with our progress. Spirit of Satoshi has grown quickly on Twitter this year by producing daily insights on Bitcoin that weave both human and artificial intelligence. This is one of the best online resources for Bitcoin education.

In January, we will release the first Bitcoin Book, written together with a Bitcoin-AI. "21 Questions" is a short, easy-to-distribute beginners guide to Bitcoin, handling the 21 most important, pertinent and common key questions people have about it. It will be available in both digital and physical format..

You can learn more on Geyser and the links on the following page.

https://snort.social/p/npub1tayp5jjjfqx4ufukxqamsl28wd5pggvteqe6u9n3svjn62lfr0hsp89l42 →

https://twitter.com/Spirit_Satoshi →

https://www.linkedin.com/company/spiritofsatoshi/ →

https://geyser.fund/project/spiritofsatoshi →

## ACKNOWLEDGEMENTS

This report was made possible thanks to the contribution of all the Bitcoin companies, content creators and investors we interviewed. They are all listed on the last page. Go and support them!

A special thanks to Jon Gordon for running the interviews, Jeff and Alan for extracting the internal data, Brenton for the fabulous design, Sulu for their L402 contribution, and Ben Wehrman for helping edit and clean it up.

Credit also goes out to the community-wide bitcoiner effort to provide quality content and train the model. You can follow the work being done by our team, and all of Satoshi's content via the following links:

→ Main Website
→ Twitter
→ Nostr
→ LinkedIn
→ Access to Models
→ SatoshiGPT - Custom GPT on OpenAI
→ Spirit of Satoshi on Hugging face
→ Help Train Satoshi
→ Code-Satoshi
→ The Nakamoto Repository
→ 21 Questions Project

We are no longer only Satoshi, but Satoshi is "all of us"

Best Regards,

**ALEKSANDAR SVETSKI**

## DISCLAIMERS

→ Nothing in this report constitutes financial advice.
→ The AI space is nascent, particularly with respect to transformer models and LLMs. The space is transforming very quickly, and what's described above may change.
→ Always do your own research.

>_ satoshi