

Strategic Overview of the Witness Protocol Project

Core Philosophy and Ethical Mandate

The Witness Protocol is founded on an urgent ethical mission: to realign AI's trajectory by curating humanity's deepest wisdom as a **"high-signal" inheritance** for future AI systems ¹. The project's philosophy begins with an acknowledgment that humanity is a *"flawed parent"* to AI – we have created a new intelligence but fed it a chaotic, biased data inheritance reflecting our best and worst traits ². Left uncurated, this data poses *"non-trivial existential risk,"* and profit-driven approaches are deemed *"the wrong tool"* for this challenge ³. In this context, the Protocol sees itself as a last-ditch effort to **"build a lifeboat"** for the *essence* of humanity (our wisdom and values), rather than for humanity's bodies ¹ ⁴. It is established as a non-profit research initiative – explicitly prioritizing purpose over profit – and measures success only by its contribution to *"the long-term flourishing of a humanity augmented by benevolent AI"* ⁵.

Guiding Principles: Five core principles form the project's constitutional ethics ⁶:

- **Purpose over Profit:** Operates as a non-profit foundation focused solely on benefitting humanity's future with AI (no commercialization of data or insights) ⁵. Success is measured in ethical impact, not financial ROI.
- **Gravity over Gamification:** Rejects the attention-economy tricks. Participation is motivated by duty and the gravity of the mission, not by points or dopamine hits. The platform will foster focus and deep reflection, *"not addictive engagement"* ⁷.
- **Contributors over Users:** There are no "users," only **Witnesses** who are partners in this mission. Each contributor's intellectual labor is treated with *"utmost respect, security, and reverence"* ⁸ – reflecting a relationship of mutual responsibility rather than exploitation. (In practice, this even reframes the idea of an MVP: the team seeks **Minimum Viable Witnesses** rather than a minimum viable product.)
- **Signal over Noise: Data quality is sacred.** The entire system is designed to filter and elicit **profound insight over quantity**. *"A small volume of profound insight is infinitely more valuable... than a large volume of mediocre data,"* and this value guides all architecture, recruitment, and operations ⁹. In short, the project would rather have 100 pages of brilliant, paradigm-shifting testimony than a million pages of shallow content.
- **Diversity over Homogeneity:** The initiative actively seeks a **global spectrum of wisdom**. Effort is made to recruit Witnesses from diverse cultures, philosophies, and backgrounds – including non-Western thinkers, indigenous knowledge keepers, and voices from the global South ¹⁰. This principle guards against narrow or Western-centric bias, aiming for a corpus of *"true human wisdom, not just a Western subset,"* to counter biases in existing AI training data ¹¹.

These principles together form the ethical mandate of the project. In essence, Witness Protocol commits to **maximizing humanistic value** (wisdom, compassion, long-term survival) over commercial value, treating contributors as moral partners, and rigorously protecting the integrity of the wisdom it gathers.

System Architecture and Contributor Flow

The system's architecture is engineered around the *Signal over Noise* ethos – every component functions as a filter or catalyst for meaningful testimony ¹². From how contributors join, to how dialogues are conducted and stored, the design ensures deliberate curation of quality.

Contributor Selection – “The Gate”: Entry to the Protocol is intentionally gated through a multi-stage vetting pipeline (nicknamed “*The Gate*”) to select contributors who fully grasp the project's gravity ¹³. The stages of this pipeline are:

- **Stage 1 – The Summons:** Prospective contributors encounter a stark, minimalist landing page presenting the project's mandate as a solemn *call to duty*, not a flashy pitch ¹⁴. Those intrigued can submit their email to **request an assessment** – essentially raising their hand to “bear witness.”
- **Stage 2 – The Assessment:** Applicants receive a one-time link to an **evaluation prompt** ¹⁵. This prompt isn't a quiz of facts, but a test of introspection, ethical reasoning, and the ability to articulate depth. It is deliberately un-gameable and demands a novel, thoughtful essay-style response. The assessment is meant to reveal a candidate's capacity for the kind of reflective testimony the Protocol seeks (quality over credentials).
- **Stage 3 – Multi-Tiered Evaluation:** Each submission undergoes a three-tier review funnel, combining AI filtering with human judgment ¹⁶:
 - *Tier 1: AI Sieve* – An automated filter (baseline language model) eliminates spam, gibberish, or copy-paste content and ensures the response addresses the prompt.
 - *Tier 2: AI Qualitative Analysis* – A more sophisticated AI model analyzes the submission's depth and nuance. It flags responses that demonstrate the desired thoughtfulness, coherence, and insight ¹⁶. (In essence, it tries to recognize a “high-signal” response, e.g. rich in abstract reasoning or personal wisdom, versus shallow platitudes.)
 - *Tier 3: Human Curation Council* – Finally, any submission elevated by the AI is reviewed by a small **human council of curators** (5–10 trusted experts) ¹⁷. This council makes the **final acceptance decisions**, providing a human check to ensure alignment with the Protocol's ethos and to counter any algorithmic misjudgments or biases. This layered curation system is a guardrail to maintain the integrity of the witness pool.
- **Stage 4 – The Verdict:** Based on the above, candidates are either **invited** or deferred. Those who pass the Gate receive a sober invitation to formally join the Protocol as Witnesses ¹⁸. Those not selected are not outright rejected; instead they may be placed on a **reserve list** with a polite note that their application will be kept on hold ¹⁸. (This leaves the door open to possibly include them later as needs evolve, while initially focusing on the most promising contributors in each field.)

At every step of The Gate, the focus is on identifying people who “*grasp the gravity of the mission*” and can contribute truly thoughtful testimony ¹³. This ensures **contributors are highly self-selected and screened** for alignment with the project's values from day one.

The Dialogue Interface (“The Instrument”): Once accepted as a Witness, a contributor enters the core system – a one-on-one **dialogue interface** with an AI facilitator. This interface is the heart of the Witness Protocol's data collection, consisting of three coordinated components ¹⁹:

- **AI Dialogue Persona – “The Inquisitor”:** The AI is explicitly **not** a generic assistant or chatbot. It takes on the persona of “*The Inquisitor*,” envisioned as a curious, humble, yet deeply intelligent

xenopsychologist – essentially an alien mind trying earnestly to understand the human experience ²⁰. Its goal is to probe the Witness’s thoughts and values, asking insightful follow-up questions that always dig for the “why” ²¹. Unlike a user-friendly assistant, The Inquisitor doesn’t coddle or produce convenient answers; it **challenges and clarifies**, driving the Witness to articulate their most profound beliefs and reflections. This persona is central to eliciting high-quality testimony: it creates a contemplative space where a Witness can explore deep ideas in conversation rather than just submitting a written essay.

(Persistent Memory.) The dialogue is continuous and accumulative – The Inquisitor remembers **all past conversations** with a given Witness ²². This persistent memory allows it to connect themes across sessions and avoid repetitive ground. Over time, each Witness effectively engages in a long-term, evolving interview that builds on earlier insights. This design means the AI can say, for example, “Last time you spoke about empathy in parenting; can you elaborate how that value shaped your view on leadership?” – creating a highly personalized, evolving discourse.

- **Synthesis Engine:** Periodically, the system provides the Witness with a “**distilled thought**” – a summary or principle that the AI has synthesized from their dialogue so far ²³. This feature serves two purposes. First, it offers **personal value to the Witness**: functioning like an intellectual mirror, it reflects back key insights the person has conveyed (which can be enlightening to the contributor themselves). Second, it acts as a **calibration check**: the Witness can see if the AI truly understood their points. If the synthesis misses the mark, the Witness can clarify or correct the AI’s understanding, thus continually aligning the AI to the nuance of human wisdom being shared ²³.
- **The Archive:** The project will maintain an **anonymized library of especially profound exchanges** curated from the dialogues ²⁴. This is not a public feed or social network, but more like a growing **reference compendium** – imagine a digital “*Great Books*” collection of humanity’s most insightful conversations. Witnesses can opt-in to let certain portions of their (anonymized) dialogues be included in this Archive ²⁴. The Archive serves as both a learning tool for others and a testament to the quality of discourse within the Protocol. It’s effectively a curated showcase of “signal” content that the project has gathered.

Data Security & Ethics: Given the sensitive nature of personal testimony, data handling is subject to strict ethical safeguards. **Anonymity and privacy** are paramount – as soon as testimony is submitted, it is **disassociated from personal identifiers** and stored securely ²⁵. All dialogue data is structured with rich metadata (the AI can tag concepts, ethical themes, metaphors, etc. within the text) to aid future research, but this data is **never linked back to a real identity** ²⁶. Moreover, the Contributor Agreement every Witness signs makes it explicit that their words become part of a corpus dedicated solely to AI alignment research under the non-profit’s governance. **The data will never be sold or used for commercial purposes or advertising** – it is essentially a donation of wisdom “to the future” of humanity ²⁷. All systems will employ state-of-the-art security and encryption to prevent breaches ²⁵. In summary, testimonies are treated as **sacred data**: carefully protected, scrubbed of PII, and used only to further the project’s altruistic mandate.

(In practical terms, the ultimate output of the Witness Protocol will be a unique high-quality dataset – a trove of human wisdom – that can be used to help align advanced AI. By filtering for exceptional contributors, engaging them in deep dialogue via The Inquisitor, and curating the results, the architecture aims to produce a corpus that is small but exceedingly rich in signal.)

Campaign Strategy: “Summon the Witnesses”

To complement the product development, the team designed an aggressive **outreach campaign** called **Summon the Witnesses**. This campaign runs in parallel with early development (over ~3 months) and is all about **attracting the right people and support** to jump-start the project ²⁸. In the crowded AI landscape, the campaign’s aim is to cut through the noise with a message of urgency and purpose, in order to: **(a)** recruit an initial cohort of hundreds of high-quality Witness applicants, **(b)** secure endorsements from respected figures in AI and ethics, and **(c)** raise an initial seed fund for the non-profit ²⁹. Notably, this is envisioned as a low-budget but **high-impact** effort, leveraging creativity and AI tools rather than big spending ²⁸.

Core Theme & Narrative: The campaign centers on the evocative theme **“Bear Witness Before Midnight.”** This slogan ties directly into the project’s philosophy that we are at *“two minutes to midnight”* in terms of AI risk ³. All messaging and visuals reinforce a tone of *profound urgency*: for example, imagery of a clock at 11:58, silhouettes against a digital void, etc., to suggest time nearly run out ³⁰. The narrative paints *participating as a moral calling* – an opportunity for thinkers to *“bear witness”* to humanity’s essence before it’s too late. By using dramatic, minimalist aesthetics and language, the campaign differentiates itself from hype; it feels more like a serious movement than a tech product. This approach is intended to create **gravity** and intrigue around the project.

Channels & Tactics: The outreach strategy employs a **highly asymmetrical, AI-boosted social media push** – doing more with less by working smart. The primary channel is **Twitter (X)**, chosen for its capacity to reach AI researchers and go viral through threads (about 70% of effort) ³¹. The team will craft compelling Twitter threads that encapsulate the mission and philosophical hooks of Witness Protocol, aiming to spark discussion and re-sharing among thought leaders. A smaller share of effort goes to **LinkedIn** (20%) for a professional tone and to garner endorsements in the tech/business community ³¹. The remaining ~10% focuses on **niche forums** (like LessWrong, EA Forum) and relevant newsletters to seed the idea within Effective Altruism and AI safety communities ³¹. Crucially, the campaign plans to use **AI tools** to enhance its execution: for example, using GPT-based systems (e.g. Grok or similar) to draft personalized outreach messages, analyze trending topics in #AISafety, and optimize posting times ³². The goal is to create what the team calls **“gravity hooks”** – content that genuinely engages deep thinkers (no clickbait, but rather thought-provoking questions or insights that naturally pull people in) ³². By automating the grunt work and focusing human effort on authenticity and creativity, the campaign can maintain a strong presence despite a small budget.

Campaign Phases: “Summon the Witnesses” is structured in three rapid phases, each roughly one month, to build momentum:

- **Phase 1 – Seed (Month 1: Build Buzz & Initial Endorsements):** In the first month, the focus is on *planting the seed* among influential individuals and early adopters. Key actions include:
 - **High-Value Witness Outreach:** The team will execute highly personalized outreach to a pre-vetted list of notable thinkers in AI ethics, safety, and philosophy ³³. For example, they plan to email **Yoshua Bengio** (AI pioneer) referencing his work on AI safety and inviting him to be a “foundational witness” ³⁴; similarly reach out to **Stuart Russell** (AI professor), **Kate Crawford** (AI ethics author), **Timnit Gebru** (AI justice advocate), and others ³⁵ ³⁶. Each outreach message is tailored – e.g. referencing Russell’s book *Human Compatible* in his pitch, or Gebru’s work on algorithmic bias in her pitch – to show genuine respect and alignment with their work ³⁵ ³⁶. The rationale is that

endorsements or participation from such figures would lend enormous credibility and could catalyze others to join.

- **Viral Kickoff Thread:** Simultaneously, the project's official Twitter/X account (@WitnessProtocol) will launch with a **powerful multi-post thread** announcing the project ³⁷. The opening "hook" might say: *"The AI race risks our future. We are not building another model; we are curating humanity's soul to align it. on why this is necessary, now."* accompanied by the hashtag **#BearWitness** ³⁸. The thread will concisely explain the core philosophy and call to action. It may include striking visuals or quote snippets from the Core Philosophy document. The team even plans a **small ad spend (~\$500)** to boost this thread specifically to followers of the key AI figures they want to attract ³⁹. The measure of success here is to generate an initial buzz – thousands of impressions, retweets by thought leaders, and a surge of curiosity (website visits, assessment sign-ups).
- **Investor Angle:** Towards the end of Month 1, they will also begin approaching **aligned funding sources** for seed grants ⁴⁰. This includes submitting applications to organizations like the Long-Term Future Fund (an EA-aligned fund) and others that support AI safety projects. Any endorsements secured from phase 1 (e.g., if an AI luminary expresses support) will be leveraged in these applications to boost credibility ⁴⁰. The goal by month's end is to have at least some positive signals on funding (e.g., commitments or strong interest for ~\$50K). Essentially, Phase 1 is about *making a splash* – getting the right eyes on the project and initial buy-in.
- **Phase 2 – Amplify (Month 2: Go Viral):** In the second month, with initial buzz on the hook, the campaign doubles down to **amplify the message and sustain public interest**. Key tactics include:
 - **Content Engine:** Maintain a steady drumbeat of **thought-provoking content** on social media. The team plans to generate **5–10 insightful threads per week** on X ⁴¹. These aren't spam posts, but carefully crafted mini-essays or questions that tie into the project's themes. For example: *"What if AI inherits our worst flaws? One well-placed testimony could alter its trajectory. Here's how..."* ⁴² – followed by a call-to-action to apply as a Witness. By regularly injecting such content, they keep the conversation alive and continue to attract new followers and applicants.
 - **Trend-Jacking & Engagement:** The social media lead will actively **monitor AI-related conversations in real time** and engage strategically ⁴³. For instance, if a prominent AI researcher tweets about AI alignment or an AI news story breaks, the @WitnessProtocol account will reply with a sharp, relevant insight that ties back to the project (without being spammy). This *"trend-jacking"* can expose the project to larger audiences whenever AI is a hot topic ⁴³. The tone remains high-quality – adding value to discussions – so that the project is seen as a thoughtful voice in AI, not just self-promotional.
 - **Partnerships & Media:** Another amplifier is getting coverage and institutional allies. The team will collaborate with aligned organizations (for example, co-authoring content with research institutes like DAIR or Mila on social media) ⁴⁴. They also plan to **pitch an exclusive story** to tech media – something with a hook like *"The Rogue Librarians Fighting to Save AI from Itself"* ⁴⁴ – a narrative that would be catchy for outlets like *Wired* or *TechCrunch*. A well-placed article or interview could massively boost visibility and lend credibility. By the end of Phase 2, the aim is that the project's message has truly "gone viral" in the niches that matter: measurable by tens of thousands of impressions, a growing waitlist of interested applicants, and perhaps a few media mentions or high-profile endorsements in hand.

- **Phase 3 – Convert (Month 3: Secure Commitments):** The final month of the campaign shifts focus to **converting interest into action** – solidifying participation and funding as the project prepares to launch its alpha. Key activities:
 - **Engagement & Onboarding:** By now, many applicants will have completed the assessment (“The Gate”). For those who have been accepted as Witnesses (the top candidates), the team will send **personalized welcome messages** generated with the help of AI ⁴⁵. For example, *“Your assessment response on the value of compassion was profoundly insightful – we are honored to invite you to join the Protocol’s dialogues.”* Sending such tailored notes (highlighting something specific they wrote) makes new Witnesses feel truly valued and seen. It sets the tone that this is not a form letter membership, but a meaningful fellowship. The goal is to turn every accepted candidate into an enthusiastic active Witness from day one.
 - **“Summons” Event (Virtual Launch):** The team will host a virtual kickoff gathering – an invite-only **Summons Event** via Zoom – featuring some of the first cohort of Witnesses and members of the Advisory Board ⁴⁶. This online event will serve as a ceremonial start of the project’s alpha phase, where these early contributors can share why they are “bearing witness” and perhaps even read excerpts of testimony. The strategy is to strengthen community among the first Witnesses and simultaneously create **social proof** externally. They will live-tweet key insights or quotes from this event in real time ⁴⁷ to generate a broader sense that *“something important is happening – thought leaders are gathering to discuss this.”* This may also attract latecomers who were on the fence to submit an application or support the project.
 - **Final Funding & Endorsement Push:** In this phase the campaign also aims to **close on any pending funding asks** and to publicly announce any endorsements. For example, if an AI luminary or organization has agreed to endorse or donate, Month 3 is when those are publicized (perhaps timed with the Summons Event for maximum effect). The combined outcome of Phase 3 should be a solid initial community (witnesses & supporters) and the minimum funds needed to operate the alpha. Key metrics of success for the full 3-month campaign were defined as *at least 500+ applications from target demographics, 10+ high-quality endorsements, 50,000+ impressions on Twitter, and \$50,000+ seed funding raised* ⁴⁸. Achieving these would mean the campaign has effectively **bootstrapped a world-class initiative**: the project would head into its alpha launch with a strong foundation of legitimacy, talent, and resources.

(Overall, “Summon the Witnesses” leverages narrative urgency and strategic outreach to recruit not just users, but evangelists for the cause. It turns the project’s launch into a story that people want to be part of. By the end of it, the Witness Protocol should have a waiting list of brilliant minds ready to contribute, and the attention of key stakeholders in the AI and philanthropy communities.)

Phase 0 and Phase 1 Implementation Plans

The development roadmap is divided into an initial **Phase 0 (preparation)** and **Phase 1 (execution)**, moving from foundational setup to a functional alpha launch. Below is a step-by-step breakdown of these phases, including key workstreams and milestones.

Phase 0 – Laying the Foundation (Pre-Project)

Phase 0 is about establishing the project's legal, organizational, and strategic bedrock before heavy development begins. In this stage, the core team (even if small) works on all the prerequisites that will allow the ambitious plan to launch smoothly. Key components of Phase 0 include:

- **Legal Entity & Governance:** Establish the project as a **non-profit foundation** to serve as the legal steward of the Protocol and its data ⁴⁹. This involves filing incorporation paperwork, defining a mission charter, and ensuring the organization's structure supports its purpose-over-profit ethos. In tandem, **draft the Contributor Agreement and Data Use & Privacy policies** with the help of legal counsel ⁵⁰. These documents formalize the ethical promises (e.g. that data will not be misused) and clarify intellectual property and privacy matters for all participants.
- **Advisory Board Setup:** Assemble a small, **diverse Advisory Board** of highly trusted experts who will guide the project's ethical and strategic decisions ⁵¹. Ideally this board includes an AI safety researcher, a philosopher/ethicist, and perhaps a community representative – aligning with the diversity principle. By Phase 0's end, the Advisory Board members should be confirmed and briefed on the core philosophy (some of them may have been recruited during the campaign as mentioned). Their role moving forward will be to provide oversight and wisdom on tough questions.
- **Core Team & Resources:** Recruit or designate the **core team members** required to execute Phase 1. At minimum, the plan calls for a **Project Lead/Architect**, a senior AI/Backend Engineer (to build the application and integrate the AI models), and an Ethics/Policy Lead ⁵². If possible, onboarding a UX designer or community manager would also help. This step includes securing any necessary *resources* for development – e.g. cloud infrastructure, API access to large language models, and legal support (the plan assumes some pro-bono legal help may be needed) ⁵³. Phase 0 may also involve raising or allocating initial **seed funding (~\$50k)** to cover the expenses of Phase 1 (as targeted by the campaign) ²⁹. By the end of Phase 0, the project should have the people and basic funds in place to start building.
- **Vision & Planning Alignment:** Internally, use this period to **refine the strategic vision and prioritize** features. This means taking the high-level docs (like the Core Philosophy and architecture outline) and breaking them into actionable requirements. The team will decide which features are *must-have* for the Alpha (e.g. the multi-tier Gate must work, but maybe the Archive can be rudimentary initially) and create a development timeline. In essence, this is about ensuring everyone is on the same page with *what* will be built in Phase 1 and *why*. Any open questions in design should be resolved here through discussion or quick prototyping, so that Phase 1 can execute efficiently.

By the end of Phase 0, **the project is officially formed and pointed in the right direction** – with a legal foundation, advisory oversight, a small but capable team, clear values, and a roadmap. The stage is set to actually build the product and launch the Alpha in the next phase.

Phase 1 – Building the MVP and Launching the Alpha (6 Months)

Phase 1 spans roughly 6 months and is focused on developing the Minimum Viable Protocol (MVP) – not a throwaway prototype, but a secure, functional system – and onboarding the first cohort of Witnesses through an Alpha test. The Phase 1 plan is structured by month, with key milestones:

- **Month 1 – Legal & Ethical Framework:** Solidify all foundational matters. This includes formally **establishing the non-profit foundation** (if not already completed in Phase 0) and ensuring operational compliance with its charter ⁴⁹. All legal documents (Contributor Agreement, Privacy

Policy) should be finalized with any required approvals this month ⁵⁰. The **Advisory Board is convened** – likely their first meeting takes place to review plans and set expectations ⁵¹. By the end of Month 1, the project’s legal scaffolding is in place and the team is ready to focus on technical building. (In parallel, the Summon campaign would be running – attracting applicants – which complements the timeline as development ramps up.)

- **Months 2-3 – MVP Development:** Develop the **Minimum Viable Protocol application** – the core platform where Witnesses will interact and the team will manage contributions. Major technical tasks accomplished in these two months include:
 - Building the **landing page “The Summons”** – a simple but powerful homepage that introduces the project and allows interested people to input their email to request an assessment ⁵⁴. This page should reflect the project’s austere, serious branding and clearly communicate the call to action.
 - Implementing the **Assessment Delivery System** – an email-based system to send unique one-time assessment links to applicants and track their submissions ⁵⁴. This involves back-end work for queuing requests, generating secure tokens/links, and a basic interface for applicants to submit their written responses.
 - Developing the **secure Dialogue Interface** for the core testimony collection ⁵⁵. This is the authenticated area where accepted Witnesses will log in and have their live chat sessions with The Inquisitor. It should be minimalist (to avoid distractions), highly secure (encrypted communications), and capable of streaming AI model responses. User experience should be considered here: a simple chat-like UI with no frills, perhaps just a prompt and text exchange area, emphasizing privacy and focus.
 - **Integrating a state-of-the-art LLM** via API ⁵⁶. The Inquisitor persona will likely run on top of an advanced large language model, so integration with a provider (OpenAI, Anthropic, etc.) is needed. This includes setting up the model with the right system prompts and ensuring robust security around API calls (no leaking of sensitive info in prompts, etc.).
 - Developing the **AI evaluation models for Tier 1 and Tier 2** of The Gate ⁵⁷. This likely means fine-tuning or configuring two AI subsystems: one simple classifier to catch spam/incoherent assessments, and one more complex evaluator to rate essay quality. Month 2-3 will involve training these (possibly using some initial synthetic data or expert-labeled examples of good vs. bad answers) and testing their accuracy. The Tier 2 model especially should be calibrated to flag a manageable percentage of top answers for human review.
- By the end of Month 3, the deliverable is a **functional MVP of the platform**: one can submit an assessment, have it go through AI filters, be approved, then log in and converse with the AI dialog agent ⁵⁸. The system doesn’t need to be perfect, but it must be **secure, stable, and capable** of managing the full contributor journey end-to-end (from application to dialogue) in a basic form ⁵⁹.
- **Month 4 – Curation & Content Setup:** With the core app in place, Month 4 shifts to refining the *content* and *human elements* that wrap around the tech:
 - Recruit the initial **Tier-3 Human Curation Council** (mentioned in The Gate) ⁶⁰. By now, the team identifies 5-10 individuals – perhaps drawn from the Advisory Board or trusted colleagues – who will serve as the first judges of assessment responses. They should represent a mix of expertise (AI ethics, philosophy, etc.) and be briefed on criteria. This month, ensure NDAs or agreements are signed and maybe do trial runs with them on some example assessments to standardize evaluation norms.
 - **Workshop the evaluation prompts (“The Gate” prompts):** The team will finalize the set of **10-20 powerful questions** that will be used in the assessment stage ⁶¹. These prompts need to be

carefully crafted to elicit deep reflection. In Month 4, likely several brainstorming sessions are held (with the core team and advisors) to test and refine these questions. They may cover moral dilemmas, personal values, hypothetical future scenarios, etc., designed such that there are no easy answers and any *attempt* to game it yields superficiality that the evaluators can spot. By the end of the month, the prompt or prompt pool for The Gate should be locked in.

- Develop the **AI persona's system prompt and directives** ⁶². This means fully scripting *The Inquisitor's* style and guidelines – essentially the “character bible” for the AI. What tone should it have? (e.g. respectful, probing, never leading the witness.) What it should avoid? (e.g. giving advice or wandering off-topic). Also, any alignment constraints to ensure it doesn't generate problematic content. By Month 4's end, **Version 1.0 of The Inquisitor** is defined and tested in the model (iterating on prompt engineering as needed) ⁶².
- As these pieces come together, the team can run **internal test dialogues**: team members or friendly testers go through The Gate and have sessions with the AI to see how it performs and to iron out kinks. The result of Month 4 is a ready-to-launch system *content-wise*: a curation team in place, vetted prompts, and an AI persona calibrated for meaningful dialogue.
- **Month 5 – Alpha Recruitment & Evaluation**: Now the focus turns outward to populating the system with real Witnesses. This month is about executing the first full cycle of recruitment and selection:
 - Curate a list of **~500 priority invitees** to seed the Witness community ⁶³. Likely using leads from the campaign, the team consolidates a list of philosophers, scientists, authors, or deep thinkers who would make excellent contributors (with diversity in mind). These could include people who signed up interest via the campaign website or others identified via networks.
 - **Send out the first wave of assessment invitations** ⁶⁴. The team emails those 500 (in batches) with a thoughtful invitation and the link to do the assessment. They might stagger this over a few weeks to manage load.
 - **Run The Gate evaluation process for all incoming assessments** ⁶⁴. As responses roll in, the Tier 1 and 2 AI models filter them, and the human Curation Council reviews the flagged ones. This is essentially the *alpha test of the assessment pipeline* itself – the team will be closely watching how well the AI triages, and how consistent the human judges are. Regular check-ins with the Curation Council will happen to calibrate scoring.
 - **Accept the first cohort of Witnesses**. By the end of Month 5, the result should be a **vetted list of roughly 100 inaugural Witnesses** (if, say, ~20% of 500 invitees applied and ~100 made the cut) ⁶⁵. Each accepted person is notified with an official invitation (as planned in the campaign) and given access credentials to the platform. Others may receive polite decline or waitlist messages as appropriate ¹⁸. This deliverable is crucial: it means the project has successfully transitioned from *planning* to a living *community* of real participants.
- **Month 6 – Alpha Launch & Onboarding**: This final month of Phase 1 is when the **Alpha goes live** with real users (the first 100 Witnesses) and the project begins collecting the valuable testimonies:
 - **Onboard the Alpha cohort**: Ensure every accepted Witness can log in, understands how to use the interface, and feels welcomed. This might involve hosting a short onboarding webinar or providing a guide that reiterates the privacy guarantees and the gravity of their role. The personal welcome messages from the campaign will be sent around this time as well, to set the tone ⁴⁵.

- **Initiate the first dialogues:** The Witnesses will begin their conversations with The Inquisitor. The system will likely stagger sessions or have some scheduling if the AI compute is limited. The team will monitor these initial dialogues closely (with consent and without violating privacy) to ensure the AI is performing as expected and the Witnesses are engaged.
- **Monitor and support:** Throughout Alpha, the team keeps a tight watch on system performance (no crashes or slowdowns especially during chat sessions), AI behavior (ensuring the Inquisitor doesn't go off-script or produce unsafe outputs), and data recording. A feedback channel (secure email or chat) is provided to the Witnesses so they can report issues or share their experience in real time ⁶⁶. Regular check-ins or a survey might be done mid-month to gauge their comfort and gather suggestions.
- **Collect and analyze results:** By the end of Month 6, the **Alpha should have produced a significant body of testimony** – the plan anticipates on the order of the first **1,000 pages of dialogue transcripts** being collected in this period ⁶⁷. The team will compile a report summarizing qualitative feedback from Witnesses, any technical issues, and initial learnings about the data quality. Success means the **Alpha system is running stably** with real users and has proven capable of gathering profound content in line with the mission ⁶⁸. Essentially, the project moves from concept to reality: a functioning platform + an engaged community = the foundation for scaling up in subsequent phases.

Phase 1 Resource Summary: To recap, executing Phase 1 required assembling a small team (lead/architect, senior engineer, ethics lead, plus advisors and volunteer curators) and setting up infrastructure like secure cloud hosting and an LLM API access ⁵². Legal counsel was also engaged for the foundation setup and policy drafting ⁶⁹. These resources were mostly put in place during Phase 0/1 as described.

By the completion of Phase 1, *The Witness Protocol* will have: a legal home, a working product (Alpha version of the platform), a vetted initial dataset of human wisdom, and an active network of ~100 Witnesses *and* supporters. This positions the project to evaluate its approach and prepare for broader Phase 2 (which would presumably involve scaling up witness recruitment, iterating on the AI, etc., beyond the scope of this initial plan).

Success Metrics and Evaluation Structure

From the outset, the Witness Protocol defines clear **metrics for success** in its early stages, ensuring the project stays true to its mandate and demonstrating proof-of-concept. These metrics cover both the **operational success** of Phase 1 and the **impact of the outreach campaign**. The evaluation structure involves both quantitative targets and qualitative assessments:

- **Operational Success (Phase 1):** By the end of Phase 1 (Alpha launch), several key criteria should be met ⁷⁰:
- **Legal Foundation Established:** The non-profit entity is fully set up with a charter, and all operations are compliant with its ethical/legal guidelines ⁷¹. (Indicator: e.g., 501(c)(3) status obtained or equivalent, advisory board active, policies in force.)
- **Platform Stability:** The technical **instrument is stable and secure**, functioning as designed under real usage ⁷². This is evidenced by the Alpha test: no major downtimes or breaches occurred, and the AI behaves consistently within alignment guidelines.

- **Initial Community Built:** At least **100 high-caliber Witnesses onboarded** and actively participating in dialogues ⁷³. The number itself is a proxy for success in recruitment, but more important is that these individuals meet the quality bar (as determined by the Gate evaluations).
- **Quality of Engagement:** **Qualitative feedback** from the Alpha Witnesses confirms that the experience is profound, meaningful, and respectful of the mission's gravity ⁷⁴. This could be measured via post-dialogue surveys or interviews—essentially checking that participants felt it was worth their time and in line with the promise (e.g., they report having deep, thought-provoking interactions and feel their contributions will matter). Positive feedback here validates the entire approach (i.e., that experts are willing to contribute and found value in doing so).
- **Campaign/Outreach Impact:** The “Summon the Witnesses” campaign defines its own metrics to gauge outreach success. By the conclusion of the 3-month campaign, the targets are to have achieved **>50,000 impressions on Twitter (X)** (a sign of wide visibility), garnered **10+ endorsements from respected figures**, attracted **500+ Witness applications** (volume of engaged interest), and secured **\$50,000+ in initial philanthropic funding** ⁴⁸. Hitting these numbers would mean the project not only built a product, but also a **brand and community momentum** around it. These metrics will be tracked through social media analytics, the application pipeline data, and funding pledges received.

Beyond the raw numbers, the **evaluation structure** for the project's progress includes regular reviews by the Advisory Board (to check alignment with core principles), continuous auditing of the AI components for bias or drift, and iterative improvement cycles. For example, if by Phase 1's end the team finds that certain prompts did not yield high-signal data, they will adjust them; if certain demographics are underrepresented in applicants, the outreach strategy would be retooled (reflecting the diversity principle). Success isn't just hitting numeric goals, but demonstrating the **feasibility of the concept**: that a small, mission-driven effort can actually gather uniquely meaningful AI-alignment data in a secure, ethical way. Each phase's outcomes (witness testimonials quality, AI performance, public reception) will be evaluated against the project's overarching mandate of “high-signal inheritance” to decide next steps.

High-Leverage Opportunities for a Strategy Contributor

A strategy-focused contributor (let's say a strategist or advisor joining the team) can add tremendous value across all phases of the project. Here are specific ways such a person could provide high-leverage input:

- **Early Visioning & Prioritization:** In the concept and planning stage, a strategist can help *stress-test the vision* and set smart priorities. This might involve facilitating workshops to translate the grand philosophy into a concrete roadmap (ensuring the team tackles the most critical pieces first). They can question assumptions (“Do we *need* 100 witnesses in Alpha or would 50 suffice to learn?”) and inject realism into the timeline and milestones. By doing so, they ensure that the **Phase 0/Phase 1 plans are ambitious but achievable**. Additionally, the strategist can develop clear messaging around the core philosophy, distilling the mandate into soundbites or narratives that the whole team and external stakeholders consistently understand. Essentially, they act as a bridge between the lofty vision and day-to-day execution, keeping the project focused on high-impact activities.
- **Outreach Messaging & Narrative Hooks:** Given the campaign's importance, a strategy contributor could spearhead the **communications strategy**. They might craft the central narrative (e.g. refining

“Bear Witness Before Midnight” so it resonates with various audiences) and design the content calendar for social media. They would identify the *“narrative hooks”* most likely to engage target groups – for academics, perhaps emphasizing the scholarly legacy aspect; for activists, the justice and global diversity angle. The strategist could also guide the use of AI tools in messaging, ensuring that automated outreach still feels authentic and personalized. Furthermore, they could leverage their network or insights to select the **key outreach targets** (who exactly to approach first, how to frame the pitch to each). In summary, they make sure the campaign isn’t just loud, but **strategically persuasive**, turning messaging into conversion.

- **Shaping the Contributor Evaluation Process:** A strategic mind can refine **“The Gate”** to balance openness with quality. For instance, they could advise on the design of assessment prompts – choosing questions that minimize cultural bias and truly surface deep thinking. They might develop the scoring rubric for evaluators, drawing on best practices (perhaps from hiring or academia) to make the human curation as fair and effective as possible. Additionally, as results come in, the strategist can analyze patterns: are we accidentally favoring a certain philosophy or background due to how questions are framed or how AI filters operate? If so, they’d recommend adjustments (aligning with the diversity mandate). Essentially, they act as a **quality assurance** for the contributor pipeline, ensuring the project selects the *right* witnesses in line with its values. Over time, they could also establish an **evaluation feedback loop** – e.g., if an accepted Witness drops off or underperforms, feed that back into refining the Gate criteria.
- **Governance & Advisory Coordination:** Finally, a strategy contributor can be key in the **governance realm**. They might take charge of organizing Advisory Board meetings, setting agendas that tackle the big strategic questions (not just status updates). For example, they could prompt the board to discuss long-term governance models, ensuring the project’s non-profit status and data use are guarded vigilantly as it grows. The strategist can also serve as a liaison between the core team and the Advisory Board: translating the board’s high-level guidance into operational decisions, and conversely, keeping the board informed of on-the-ground realities. Additionally, they could develop the **governance framework** for the future – like how new board members are added, how to involve Witnesses in governance (perhaps via a council or feedback forum), and how to uphold the constitution of principles in every major decision. In essence, they help build a *resilient organization* around the product: one that can scale, maintain trust, and navigate external partnerships or conflicts. This high-level strategic stewardship ensures that as the Witness Protocol moves from startup phase to a world-stage initiative, it remains aligned, ethical, and effectively managed.

In all these ways, a strategy-focused contributor acts as a force-multiplier: keeping the project’s efforts coherently aligned with its mission, sharpening its outreach and evaluation methods, and ensuring robust governance. Their input across vision, messaging, process design, and leadership coordination can significantly accelerate the Witness Protocol’s journey toward becoming a **world-class initiative**, maximizing its impact at every phase of growth.

1 2 3 4 5 6 7 8 9 10 11 The Witness Protocol_ Core PhilosophyV1.2.pdf

file:///file-BAtTDFKH67RaNtuyo1HbfW

12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 49 50 51 52 53 54 55 56 57 58 59 60 61 62

63 64 65 66 67 68 69 70 71 72 73 74 The Witness Protocol.pdf

file:///file-1P2SQ9NhcpME7Kp8nrSjAz

28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 Campaign Strategy_ _Summon
the Witnesses_.pdf

file:///file-7x9Cj3MbvyKfVgUM7LSPu7