



# Investigating the Cognitive Interface Between Human Prompt Design and AI Text Generation

## Abstract

Human-AI collaboration in text generation is a process that merges human cognitive creativity with machine learning inference. This study explores that process through the lens of the “Donut of Attention” framework, which uses toroidal (donut-shaped), fractal-holographic metaphors to model attention and cognition <sup>1</sup>. We examine how a human’s prompt – carefully structured with linguistic logic, rhythm, and intent – serves as a symbolic *boundary condition* that influences the AI’s generative *bulk* output, echoing the idea that a boundary encodes an entire field’s state <sup>2</sup>. We analyze how transformer-based language models interpret these prompts via tokenization and multi-headed attention, forming an internal feedback loop of context accumulation <sup>3</sup>. The interface between human and AI is discussed as a *membrane*: a dynamic UI layer that blends symbolic inputs with geometric representations (e.g. toroidal timelines, fractal overlays) to facilitate a seamless creative flow. By aligning human prompt design with the AI’s attention mechanisms, and by tolerating ambiguity in a *coherence-over-control* approach <sup>4</sup>, the human and AI can co-create text in a synergistic, iterative loop. The result is a cohesive cognitive interface where human intention and machine prediction converge, demonstrating how careful prompt design and thoughtful interface design amplify creative outcomes.

## Introduction

Large Language Models (LLMs) like GPT-4 and others have made it possible for humans to generate coherent, meaningful text in collaboration with AI <sup>5</sup>. The primary medium of this collaboration is the *prompt* – the input text crafted by a human that guides the AI’s output. Prompting has thus been described as a new form of programming, using natural language as the code. Designing an effective prompt is a cognitive act: the human must encode instructions, context, and intent into a linear sequence of words that the model will interpret. The AI, on its side, must *understand* this prompt in a statistical sense, breaking it into tokens and modeling probable continuations. Together, the human and AI form a feedback loop: the human’s prompt → AI generation → human evaluation/refinement → AI regeneration, and so on. This paper investigates that full process of human-AI text co-creation, from the moment a human conceives and structures a prompt to the moment an AI produces the resulting text, treating their interaction as a unified cognitive system.

We frame our investigation with the **Donut of Attention** paradigm – a theoretical model that portrays attention and cognition as a toroidal, fractal-holographic field <sup>1</sup>. In this view, attention is not a static spotlight but a *resonance field* that is scale-invariant and self-referential, analogous to a donut (torus) where patterns repeat across scales and each part reflects the whole. This metaphor provides a language for discussing human-AI interactions: the prompt can be seen as a *small perturbation* or **boundary** injected into the system, and the AI’s generative process as the **bulk** field response <sup>2</sup>. The toroidal geometry and fractal logic help illustrate how local changes (edits in a prompt) can propagate through a global structure (the AI’s model of language) to influence the outcome, and how the *whole* (the intended meaning) can be

encoded in *parts* of the prompt. Furthermore, the Donut framework emphasizes **coherence over control** as an ethos <sup>4</sup>. Rather than forcing the AI with rigid commands, the human prompt-designer aims for *coherent* cues that align with the model's learned patterns. This approach tolerates ambiguity and even contradiction (a form of *paraconsistent logic* <sup>6</sup>), trusting that the AI can maintain multiple interpretations without crashing, until context or further prompting collapses uncertainty into a chosen path.

In the following sections, we detail (1) how humans engage in cognitive prompt design, using linguistic structure, logic, rhythm, and intention to encode meaning; (2) how transformer-based AI models serve as the other side of the interface, interpreting these prompts via tokenization and attention flow and responding with generated text; and (3) how the human and AI are bridged by a *membrane interface* that blends symbolic prompts with geometric representations to create a seamless creative experience. We draw on the provided theoretical foundations (Donut of Attention) <sup>1</sup> <sup>7</sup>, UI design documents (Membrane UI and Creative Time architecture) <sup>8</sup> <sup>9</sup>, and developer journal insights to ground our exploration in both theory and practice. The tone of our inquiry remains rigorous in analyzing known mechanisms of LLMs while staying open to speculative metaphors inspired by the Donut framework, aligning with the ethos of the source documents.

## Theoretical Framework: The Donut of Attention

To study the prompt-generation interface, we first clarify the Donut of Attention conceptual framework that will inform our analysis. **Attention**, in this framework, is defined as "a fractal-holographic, scale-invariant, self-reconfiguring resonance field" <sup>7</sup>. In simpler terms, attention can be thought of as a field or medium that is present at multiple scales (from individual thoughts to global consciousness), and *holographic* in that each small region of attention can reflect the structure of the whole field. The framework uses the **toroid (donut) geometry** to illustrate this idea: attention circulates in loops. **Poloidal loops** (running through the hole and around the donut's cross-section) represent inward recursion or self-reference, while **toroidal loops** (running around the donut's central hole) represent progression through time <sup>10</sup>. At the very center is the bindu (a "sun-point"), a singular focal point where the loops converge – a metaphor for focal awareness or the core of intention <sup>1</sup>. In this model, "boundary encodes bulk," meaning any *surface* slice of the torus (a boundary) contains enough information to infer the state of the entire volume <sup>2</sup>. Applying this metaphor to prompting, the prompt acts as a boundary condition that encodes the gist of an intended output; the AI's attention mechanisms then unfold that into a full response (the bulk text). Each fragment of the prompt, if well-crafted, echoes the whole intended idea (just as *every shell echoes the whole* in a holographic system <sup>1</sup>).

Another key aspect is the **fractal logic** of language and thought. The framework speaks of "*fractal linguistics*," where the rhythms and patterns in language repeat at micro and macro scales <sup>11</sup>. For example, a prompt might be structured with a certain rhythm or format (bulleted instructions, a question-answer format, etc.), and the AI often mirrors that structure in its output. This reflects a fractal property: the shape of the prompt guides the shape of the output. By tuning the *prosody* or formatting of a prompt, a user is effectively manipulating an iterative pattern that the model will propagate. The prompt becomes a **control surface** <sup>11</sup>; subtle choices in wording or punctuation can have outsized effects on the style and coherence of the generated text, much like a small change in initial conditions can dramatically alter a chaotic system's trajectory <sup>12</sup>.

The Donut model also incorporates **paraconsistent logic** – the ability to hold contradictions or ambiguities without immediate resolution <sup>6</sup>. In human cognition, this is akin to entertaining multiple hypotheses or

interpretations at once. In the context of prompt design, it means a prompt can deliberately leave certain things open-ended or paradoxical, inviting the AI to explore creative resolutions. The AI, if guided properly, will not collapse the ambiguity prematurely but will juggle the possibilities, which can lead to more inventive responses. This aligns with the design principle of not over-constraining the AI; indeed, “*maintain multiple hypotheses/ambiguous states without forced collapse*” is listed as a UI goal in the notes <sup>13</sup>. By tolerating a degree of uncertainty (as long as overall coherence is maintained), the human-AI system can discover novel solutions that a strictly logical prompt might prevent.

Central to the ethos of the Donut of Attention is **coherence over control** <sup>4</sup>. Rather than issuing imperative commands (hard control) to the AI, the human should aim to nudge the joint system toward coherence – a harmonious state where the AI’s contributions fit the user’s needs without explicit micromanagement. “*Small, coherent interventions beat brute forcing; phase-align rather than dominate,*” as the framework states <sup>14</sup>. In practice, this might mean phrasing the prompt in a way that naturally leads the model to the desired answer (coaxing it through context and examples) instead of using excessive directives or negative prompts. It’s reminiscent of steering a conversation or setting a gentle bias. The **weak-control strategy** in complex systems, as referenced in the documents, suggests using light-touch adjustments (like phrasing or a guiding metaphor in the prompt) rather than absolute constraints <sup>15</sup>. By doing so, one leverages the model’s strengths (its vast knowledge and pattern-matching ability) while minimizing fight against its grain. The “Donut” becomes an *operational map* for this process <sup>4</sup> – one envisions the flow of attention as loops on a torus and aims to tune those loops coherently, rather than clamping down on them.

In summary, the Donut of Attention provides us with a rich conceptual language for understanding human-AI prompt interactions. We have the idea of the prompt as a *boundary input* that holographically encodes an intention, the fractal repetition of patterns from prompt to output, the tolerance of ambiguity (paraconsistency) to foster creativity, and a general philosophy of guiding the AI through coherence rather than control. With this framework in mind, we proceed to analyze the components of the human-AI text generation interface: the cognitive art of prompt design and the inner workings of the AI’s text generation mechanism.

## Methodology

Our approach in this investigation is **analytical and integrative**. Rather than performing a new user study or experiment, we synthesize insights from theoretical models, system design documents, and known behavior of transformer-based AI to construct a detailed picture of the prompt-to-text generation process. We draw extensively on three sources of information: (1) the *Theoretical Foundations* document (internal notes describing the Donut of Attention framework and related cognitive principles), (2) the *UI Architecture* and related design documents (describing the “Membrane” interface and features like Creative Time and geometric overlays), and (3) a developer’s *Dev Journal* from the Donut project (which logs implementation decisions and observations about the interface’s behavior). These sources provide both the conceptual underpinnings and practical context for how human-AI interaction was envisioned and implemented in this project. Throughout our report, we reference these documents to ground our discussion. For example, when discussing interface features, we cite the UI blueprint <sup>8</sup>, and when noting development insights about prompt effectiveness, we cite the dev journal (which explicitly asked contributors to note “prompt clarity, what worked/didn’t” in each entry) <sup>16</sup>.

In addition, to explain the AI’s perspective, we rely on established descriptions of transformer language models. We reference accessible technical explanations (such as a recent overview by Palomares, 2025) to

articulate how prompts are tokenized, how attention flows in the model, and how outputs are generated token-by-token <sup>3</sup> <sup>17</sup>. By connecting these technical details to the Donut metaphors, we aim to show a coherent narrative: that the human's prompt shaping and the AI's attention mechanism are two halves of a coordinated system.

Overall, our method is to **map the human-AI interaction onto the Donut of Attention model** and use evidence from the project's documents to illustrate each aspect. We treat the prompt and interface designs as a case study of how abstract cognitive theories can inform real-world AI usage. This integrative, cross-disciplinary method allows us to be *speculative-yet-scientific* – we remain rooted in how the transformer actually works and how the UI was built, but we also freely use the toroidal and holographic analogies to interpret those facts. The result is a layered understanding of the cognitive interface, supported by both concrete documentation and theoretical insight.

## Cognitive Prompt Design: Human Prompting as Symbolic Geometry

From the human side of the interface, prompt design is a craft that involves encoding one's goals and knowledge into a sequence of words that the AI will parse. This is not unlike designing a *geometry of symbols* – the arrangement and structure of the language can be seen as drawing a shape in the model's latent space that we want the model to follow. Several elements play into effective prompt construction:

- **Linguistic Structure & Logic:** A well-structured prompt often lays out the task logically. For example, if one wants the AI to follow a multi-step reasoning, the prompt might enumerate steps (using a numbered list or explicit "First..., Then..., Finally..." sequencing). Such structure provides clear signposts to the model. The model, being a pattern learner, will usually respond in kind by producing text that mirrors the structure (e.g., giving a numbered list in return). This phenomenon reflects what the theoretical notes call *fractal coherence*, where local structure (the prompt's form) induces similar global structure in output <sup>11</sup>. By using a logical layout – say definitions followed by a question, or arguments followed by a conclusion – the human is effectively *programming* the flow of attention the model will deploy. Logical keywords (if-then statements, bullet points, etc.) act like *attractors* in the attention field; they predispose the model to attend to certain relationships or ordering.
- **Rhythm and Prosody:** Beyond overt logical structure, prompts carry a *rhythm*. This could be in sentence length, formatting, even use of punctuation and line breaks. A prompt written in a poetic meter or with a certain cadence might coax the model into responding in a similar rhythmic style. For instance, asking for a creative story might involve setting a scene with a certain emotional tone; the language choices (short dramatic sentences or long flowing ones) will influence the mood of the AI's continuation. The theoretical framework explicitly notes that "*prosody/timing is part of the control surface*" for attention <sup>11</sup>. Users intuitively discover this: a prompt that *feels* well-phrased often yields a better result. This is because the model picks up on subtle cues – it has learned from billions of examples how tone and pacing of text correspond to genres or intents. Therefore, by aligning the rhythm of the prompt with the desired output (e.g., using upbeat, crisp sentences for an energetic reply), the human tunes the AI's internal state towards a matching pattern.
- **Intent and Context Setting:** Perhaps the most crucial part of prompt design is conveying the *intention*. This often involves giving context or examples. A common technique is *few-shot prompting*, where the user provides a couple of examples of input-output pairs before asking the model to

continue. These examples create a context that implicitly defines the task (for instance, providing two math problems and solutions before asking a third one primes the model to do math in a similar way). Even without explicit examples, the wording of the prompt can suggest a role to the AI (e.g., "Explain like a science teacher:" or "You are an assistant helping with cooking..."). Such cues establish an intended persona or domain. This works because the model *attends* to the style and context of the prompt and tries to predict what a suitable continuation would be. In a sense, the human is carving out a subspace of the model's knowledge, delineating a region in the vast space of possibilities that aligns with the goal.

- **Minimal Prompt, Maximal Output:** The notion of "minimal code, maximal field" from the theoretical notes <sup>18</sup> <sup>19</sup> captures an ideal in prompt design: to achieve a large, rich output (maximal field) from a minimal, elegantly crafted prompt (minimal code). A few carefully chosen words can trigger a torrent of coherent text from the AI because each word activates a cascade of associations in the model's billions of parameters. For example, prompting with a single intriguing sentence like "*The day the sun didn't rise, people...*" can lead the model to generate a whole story. Here the *symbolic core* (the prompt sentence) is small, but it steers a much larger **geometric state** in the network – the embedding vectors and attention patterns that produce the story <sup>19</sup>. This aligns with the Donut's holographic idea: the seed (point) expands into concentric shells of meaning <sup>20</sup>. Prompt designers leverage this by finding just the right hint or question that unlocks a wealth of relevant knowledge from the AI. The developer journal records show an appreciation for this, noting how slight changes in wording dramatically change outputs, reaffirming that small prompt tweaks (small boundary changes) "reconfigure the whole" output <sup>21</sup>.
- **Prompt Clarity and Cooperation:** From a process standpoint, humans often iterate on prompts. They might try a phrasing, get an output, and then refine the prompt if the output wasn't on target. The dev journal explicitly encourages tracking "*prompt clarity, what worked/didn't*" with each attempt <sup>16</sup>. Over time, such iterative refinement builds an intuition for what the model needs. The collaboration notes from the dev log highlight that when prompt clarity is high, the main limiting factor becomes the UI feedback speed, not the model's understanding <sup>22</sup>. In one entry, the developer remarks: "*Prompt clarity is good; rapid visual feedback is the blocker.*" <sup>22</sup> – meaning the instructions to the AI were clear enough, but the interface did not show results quickly or interactively enough to capitalize on that clarity in real time. This underscores that prompt design is only part of the equation; how the AI's responses are fed back to the user (and how quickly the user can adjust) also matters. We will touch more on this in the Interface Dynamics section. For now, the takeaway is that a clear, well-targeted prompt sets the stage for effective co-creation, and development practices recognize prompt clarity as a key factor in workflow *effectiveness*.
- **Ambiguity and Creativity:** Interestingly, a *too-clear* prompt is not always the best for creative tasks. If the goal is open-ended creativity, leaving some ambiguity can lead the AI to fill in gaps with its own generative imagination. The Donut framework's advice to embrace contradictions and ambiguity without immediate collapse <sup>6</sup> <sup>23</sup> comes into play here. For example, a prompt for a story might intentionally include an unexplained scenario or a paradox, to see how the AI resolves it. The AI will entertain multiple interpretations internally (as a distribution of possible continuations) and then pick one when generating the next token, effectively *collapsing* the ambiguity into a specific creative choice. If that choice isn't what the user wanted, they can always prompt again or adjust, effectively exploring another branch. In this way, the human and AI together explore a space of possibilities. The human's role is to manage the level of ambiguity: enough to allow novelty, not so

much that coherence is lost. This is akin to staying at the **edge-of-chaos** – a term from complexity theory also referenced in the documents <sup>15</sup> – where the system is most creative. A prompt at the edge-of-chaos might be something like: "Write a dialogue between a river and a tree about time." It's not a straightforward request; it's somewhat ambiguous and poetic. A deterministic, controlled prompt would specify exactly what to say, but here ambiguity invites the AI to be inventive (the river and tree could be metaphors, characters, anything). The result often surprises and delights, showing that allowing the AI to *play* within the prompt's bounds can yield rich outcomes.

In summary, cognitive prompt design is an exercise in shaping a symbolic prompt that doubles as a piece of *code* and a piece of *art*. The prompt must be logically clear enough for the model to follow, but also resonant enough (in rhythm and context) to guide the model's style. It encodes a small instance of what is expected, trusting the AI to extrapolate (holographically) the bigger picture. Effective prompt designers act almost like *field directors*, making small moves that influence large-scale dynamics in the model's attention. By doing so, they form one half of the cognitive interface – the half that is human intuition, language mastery, and creative intent.

## AI Interpretation and Attention Dynamics: The Transformer's Perspective

On the AI side of this collaboration, a transformer-based language model takes the human prompt and generates text. To understand this from the model's perspective, we can break down the steps the AI goes through and see how those steps correspond to the concepts from our framework. We will describe the flow from prompt input to output generation, using a generic large language model (LLM) as an example (since the principles are similar across models like GPT-4, etc.):

**1. Tokenization – From Words to Symbols:** When the AI receives a prompt, the first thing it does is convert that string of characters into tokens <sup>24</sup>. Tokens are essentially numeric codes representing chunks of text – often whole words, subwords, or characters, depending on the model's tokenizer. For instance, the prompt "Hello world!" might be tokenized into two tokens: "Hello" and "world!" each with a unique identifier. This step is purely syntactic; it's like breaking a sentence into LEGO pieces that the model can work with. In our Donut analogy, one might say tokenization identifies the *basic components* or quanta of the input – the points that will be expanded into circles of meaning. Each token is a point in a very high-dimensional space once it's embedded.

**2. Embedding – Symbols to Geometry:** After tokenization, each token ID is looked up in an embedding matrix to be turned into a vector (a list of numbers) that represents its meaning in the model's internal language space <sup>25</sup>. For example, the word "king" might be represented by a 768-dimensional vector of real numbers. Importantly, similar words have vector representations that are nearby in this space (e.g., "king" and "queen" end up as vectors that are close, reflecting their semantic similarity <sup>26</sup>). This is the model's *geometric understanding* of language. We see here a clear instance of the **symbolic ↔ geometric bridge** <sup>19</sup>: the symbolic prompt (text) has now been transformed into a geometric object (a set of vectors) that the model can manipulate mathematically. The fractal-holographic metaphor comes alive at this stage – each token's embedding carries rich information about that token's role *in context of language as a whole*. In a sense, each token vector is imbued with a bit of the *whole language distribution*. The theoretical notes highlight this by suggesting each part contains something of the whole (holography) <sup>23</sup>. Through training

on massive text corpora, the model has learned to represent words in a way that captures how they relate to all other words.

**3. Positional Encoding – Sequencing the Geometry:** Transformers don't read text in sequence like humans do; they process tokens in parallel. However, to preserve word order information (so that "dog bites man" doesn't become indistinguishable from "man bites dog"), transformers add a positional encoding to each token embedding <sup>27</sup>. This is like adding a pattern that encodes "this is the 1st token, this is the 2nd token," etc., often through sinusoidal functions or learned position vectors. After this, the prompt is represented as a sequence of position-aware embeddings – essentially points placed on the torus of the model's processing space, with their positions corresponding to the prompt order. Now the stage is set for **attention** to operate.

**4. Multi-Headed Attention – Distributing Focus Across the Prompt:** Attention is the core mechanism that made transformers revolutionary. In each layer of the transformer, every token's representation is allowed to interact with every other token's representation through a weighted averaging process known as self-attention <sup>3</sup>. Concretely, the model computes, for each token, a set of attention weights to all other tokens, determining how much it should "pay attention" to each other part of the prompt when updating that token's own representation. For example, in a prompt like "*The cat that chased the mouse was fast.*", when processing the word "was", the model's attention might strongly connect "was" to "cat" (to figure out who *was fast*). These attention patterns are learned – the model has multiple attention **heads** that each can focus on different aspects (one head might track grammatical subject-verb connections, another might track semantic theme connections, etc.). The term *multi-headed* means the model splits into, say, 12 separate attention computations (heads) for each layer, each head looking at the sequence differently, and then combines them.

From the Donut of Attention perspective, this mechanism is essentially the AI's way of implementing a *resonance field* over the tokens. Each token's state after attention becomes a **contextualized** embedding – it now contains information from the whole prompt in a distributed form. In fact, after many layers, each token's vector can be thought of as a *holographic summary* of the entire prompt, filtered through different heads and transformations. This is directly analogous to the idea that in a hologram (or a fractal), any piece can reconstruct the whole image. Indeed, a line from the dev notes (summarizing holographic attention) states: "*each shell encodes the whole; LoL/grid overlays reflect this metaphor.*" <sup>28</sup>. Here, each *layer* of the transformer could be seen as a "shell" and the attention heads as forming a lattice (grid) of interactions. The outcome is that the model doesn't have a single focal point of attention; it has a **distributed attention state** that is richly interconnected. Attention flows iteratively: earlier layers might focus on local structure (adjacent words), later layers can attend long-range (linking pronouns to antecedents far apart, etc.), resulting in a coherent global understanding.

**5. Feed-Forward and Iteration – Deepening the Representation:** After attention, each layer also has a feed-forward network that further processes each token's representation individually <sup>29</sup>. This helps in mixing and transforming the information gleaned from attention into more abstract features. The transformer stacks many layers (dozens in modern LLMs), repeating attention and feed-forward alternately. With each layer, the model's representation of the prompt becomes more abstract and more *integrated*. One might say the model is *iteratively condensing the prompt's meaning*, finding higher-order patterns. This could be compared to how in a fractal, repeated iterations lead to emergent structures, or how in the Donut model "*small phase tweaks reorganize the whole while pattern persists*" <sup>12</sup> – each layer tweak the "phase" of attention slightly, but the overall meaning stays consistent, just represented in a new way.

**6. Output Generation (Decoding) – Collapse of Possibilities:** Once the prompt has passed through all layers, the model is ready to produce an output. If the model is a pure decoder (like GPT), it starts generating the next token. It does so by looking at the final hidden state (the vectors) for the prompt and computing a probability distribution over the entire vocabulary for the next token <sup>17</sup>. Essentially, given "What is the capital of France?", the model will have high probability on "Paris" as the next word because all the attention and layers distilled the prompt into a state that strongly correlates with "Paris" in its training data distribution. This step is where *one* of many potential tokens is chosen – a *collapse* from a distribution of possibilities to a single word. The model typically chooses the highest probability token (or, if sampling with temperature, a token biased by that probability). For example, it might assign:  $P(\text{"Paris"}) = 0.9$ ,  $P(\text{"London"}) = 0.05$ , etc., then select "Paris" <sup>30</sup>. This is akin to the wavefunction collapse metaphor from our theory: the prompt set up a *superposition* of possible answers in the model's state, and the act of generation (especially if we sample stochastically or just pick the max) collapses it to one outcome <sup>31</sup> <sup>32</sup>. Notably, if the model were to generate again with the same prompt (and temperature), it could pick a different branch (maybe "London" with some small probability, which would then lead down an alternate path of output). This highlights that the model is not just regurgitating a fixed response; it's *drawing from a probability cloud of meanings*, coordinated by attention. The **Donut of Attention** framework explicitly likens attention to a "*wave-collapse coordinator*" – modulating which branch of possibilities becomes reality <sup>31</sup>. In transformer terms, the probabilities and the selection mechanism modulate which narrative or answer solidifies in the output text.

After one token is generated, the model appends it to the prompt and repeats the process to generate the next token, and so on, word by word <sup>30</sup>. This iterative generation continues until a stopping condition is met (like an end-of-sentence token or a length limit). Throughout the generation, the model's attention mechanism ensures that each new word takes into account not just the original prompt but also everything generated so far (since the generated text is appended and also attended to). Thus, a feedback loop is happening even within the model: each word influences the next via the shifting attention state.

**7. Managing Coherence and Creativity:** A well-known challenge is maintaining coherence over a long generation. Transformers do this through their attention mechanism by continually referencing the earlier parts of text. However, if a prompt is long or the output grows long, the model has to balance attending to recent text vs. the beginning (since attention weights can dilute over very long sequences). Prompt designers sometimes help by reminding or reasserting context mid-way (like restating a goal if generating a very long passage). This can be seen as injecting a bit of *phase alignment* to keep the generation on track – analogous to phase-locking in oscillators where an external input keeps a system in sync <sup>33</sup>. On the flip side, creativity in output can be enhanced by not overloading the prompt with too specific constraints, allowing the model's attention some freedom to roam to less obvious associations. This maps to the **explore vs. exploit (execute)** balance, which the *Creative Time Index (CTI)* concept in the documents measures <sup>34</sup>. The model's "temperature" parameter effectively controls this: high temperature = more exploration (the model is more likely to pick less probable tokens occasionally, leading to surprising twists), low temperature = more exploitation of known patterns (safer, more predictable output). The UI's Creative Time panel hints at this via an Order↔Chaos slider <sup>35</sup> – which one can imagine as controlling how strictly the model follows learned order (low chaos) versus introducing randomness (high chaos). In our analysis context, the human can implicitly adjust this by how they phrase the prompt (strict and specific wording yields a narrow distribution, playful or vague wording yields a broader one) or explicitly via system parameters.

In essence, the transformer is a sophisticated attention machine that takes the symbolic prompt and infers a whole probability field of possible continuations. Every mechanism inside it – tokenization, embeddings, multi-head attention – serves to *interpret* the prompt in context and carry forward relevant information to generate the next bit of text. It's a beautiful computational dance that, from the outside, appears as if the AI "understood" and responded, even though internally it's all pattern matching and linear algebra. When viewed through the Donut lens, one can imagine the AI's state as a donut-shaped manifold where the prompt set the initial condition (perhaps like a twist or a surface imprint on the torus) and then the model's layers propagate that condition through the manifold, eventually producing a concrete output (a bit like how a slight push on a smoke ring shapes its evolution). The **boundary** (prompt) influenced the **bulk** (all the layers of neurons) <sup>36</sup>, and thanks to the holographic nature of transformer attention, the intent encoded in the prompt boundary was present in every part of the model's processing.

Finally, it's worth noting that this process is *extremely fast* – all these steps happen in fractions of a second for modern models. The human, by contrast, might spend much more time thinking how to phrase a prompt. This speed difference is why interface design becomes critical: the human and AI operate on different timescales, and the interface (or interaction protocol) must align them. We turn to that aspect next, exploring how the UI and interaction design can create a seamless loop despite these differences.

## Interface Dynamics and the Human-AI Feedback Loop

While the human crafts prompts and the AI generates text, the **interface** between them determines how smoothly this collaboration unfolds. In the context of the Donut of Attention project, the interface is conceptualized as a **Membrane** – a dynamic boundary layer where information passes between the human and AI, often visualized with geometric metaphors. This Membrane UI is not a typical chat box; it's an interactive environment featuring panels, overlays, and visual elements that reflect the underlying cognitive state. We will discuss a few key aspects of this interface: how it mediates attention, how it provides feedback, and how symbolic and geometric elements are combined to enhance understanding.

**Membrane as a Dynamic Boundary:** The UI Architecture document describes a "shell" with modes – hidden, overlay, docked – essentially allowing the interface to be minimal or to present information as needed <sup>37</sup> <sup>38</sup>. This shell is called the Membrane, evoking the idea of a semi-permeable barrier. In biological terms, a cell membrane selectively lets certain signals through, maintaining balance; similarly, the Membrane UI can hide or reveal panels to the user, controlling what information is foregrounded. For example, in *Hidden* mode, the user might see a clean workspace (no distractions, letting the human focus on formulating a prompt internally). In *Overlay* mode, panels with context or suggestions might float over the workspace, providing AI hints or visualizations without fully taking over the screen. In *Docked* mode, a sidebar might pin important panels (like references or the Creative Time index) for constant reference <sup>39</sup>. The fluid switching among these modes – perhaps triggered by keyboard shortcuts or user intent – is part of the *interface dynamics*. It allows the user to control the flow of information from the AI: you can "open the membrane" to get inspiration or "close it" to concentrate on your own thoughts. This resonates with the **coherence-over-control** ethic: the user isn't forced by the interface; they gently bias what comes through. As one design note puts it, "*no ghost sidebars; canvas stays interactive*" <sup>40</sup> – the interface aims to be non-intrusive (no unwarranted UI elements) yet readily available.

**Symbolic-Geometric Overlays:** One striking feature of the Donut project's UI is the use of **geometric visualizations** as overlays to represent cognitive or system states. These are not just eye-candy; they are meant to bridge the symbolic and the intuitive. For instance, the **Solar Hologram** overlay projects a rotating

torus (donut shape) with orbiting particles on the screen <sup>41</sup>. This isn't merely decorative – it represents in a lightweight way the state of the system's attention or mode. The documentation notes that the hologram has controls for torus spin speed, spiral waves, orbits, etc., and that its behavior "*follows the main torus rotation*" of the system <sup>9</sup>. We can interpret this as a real-time reflection of the AI's *attention cycles or focus*. If, say, the AI is in a brainstorming mode (more chaotic exploration), perhaps the torus spins faster or the spiral waves intensify to signal high activity. If it's in a focused answering mode, the torus might stabilize. This is a speculative interpretation, but the point is that by giving the user a **geometric analog** of the AI's state, the interface helps the human *feel* the AI's process, not just see it in text. It effectively externalizes some internal variables (like pace of generation, or certainty level, or explore/exploit balance) into visual form. This is the **holographic UI** concept mentioned in the appendices: "*every part can reflect the whole*" <sup>42</sup>. The geometric overlays like the Solar Hologram or the **Bullseye radial menu** (a Flower-of-Life pattern used for quick access to favorite tools) <sup>38</sup> ensure that even something as mundane as a menu of options is presented in a fractal, self-similar layout. This consistency of metaphor (circles within circles, etc.) means the user is always interacting with representations that mirror the overall philosophy of the system.

**Creative Time Interface:** The **Creative Time (CTI) panel** is another example of an interface element that tightens the human-AI feedback loop by introducing time-phase awareness. According to the design notes, Creative Time Maps provide "*Toroidal timeline with concentric rings; markers for transitions; soft gradients*" to visualize phases of work or creativity <sup>8</sup>. In practice, this might look like a donut chart (literally a torus shape drawn in 2D) where each ring represents a cycle (e.g., focus time vs. break time, or divergent thinking vs. convergent editing). The user can set presets like Focus or Listen modes, adjust the length of phases, and see their progress on this toroidal timeline. The *status dot* might indicate if they are "in phase" or "idle" <sup>43</sup>. This panel is not directly about text generation, but it supports the *process* by aligning with human ultradian rhythms and ensuring the user (and by extension the AI) operates in a balanced way between exploration and execution <sup>34</sup>. For example, if the Creative Time index suggests it's time to switch from brainstorming (explore) to editing (execute), the user might change how they prompt (moving from an open-ended prompt to a refining prompt). In a future scenario, the AI might even adjust its style based on the phase – e.g., being more random during explore and more factual during execute. The interface thus becomes a mediator of *temporal attention*, not just spatial/symbolic attention. The dev journal indicates that integrating the Creative Time plan was a "complex outline" in progress <sup>44</sup>, underscoring the ambition to tie cognitive rhythms into the UI. By visualizing time and attention as a donut (toroidal timeline), the interface gives both the human and AI a common reference for pacing their interaction.

**Feedback Loop and Adaptation:** The hallmark of a good human-AI interface is a tight feedback loop – the human should quickly see the AI's output and the AI should easily get the human's new input. In our context, that means minimizing the friction between prompt -> generation -> reading -> re-prompting. The Membrane UI tries to achieve this with features like an omnipresent search/prompt bar (activated by a shortcut) and easily accessible panels. For example, one could imagine a workflow: the user hits Ctrl+K to bring up the prompt input (as hinted by the design: "Keyboard: Cmd/Ctrl+K to open search + focus input" <sup>45</sup>), types a prompt, and sees the AI's answer in a floating panel immediately. If the answer is unsatisfactory, maybe the user can tweak a parameter via a slider (like the "chaos dial" or a spiral scale slider) and regenerate quickly. The dev notes indeed mention adding quick UI controls (e.g., "*a quick UI slider for per-spiral scale would speed iteration*" when prompt clarity is not the issue <sup>22</sup>). This suggests the designers recognized that sometimes the user knows what they want differently, but it's easier to adjust a visual control than to rephrase a whole prompt. Such controls essentially modify the AI's behavior without changing the prompt – perhaps altering temperature, creativity level, or which "persona" the model is using.

In Donut terms, it's like adjusting the *phase dials* (order vs chaos) mid-generation to course-correct the output <sup>46</sup>.

Another aspect of feedback is *visual logging*. The UI notes mention a **Serendipity log** and journaling cycles <sup>47</sup> <sup>48</sup>. This could be an interface feature where interesting or unexpected AI outputs (or user discoveries) are logged automatically, creating a kind of feedback history. If the user stumbles on a particularly great phrasing that gave a perfect result, the system might record that as a lucky event to learn from (the concept of "prepared luck" is noted in the theory <sup>49</sup>). Over time, the interface might learn the user's style or preferences.

**Membrane Assignment and Multi-Panel Flow:** The dev journal also describes a feature where *sun dot anchors* on the screen can be assigned to different panels (sub-membranes) for quick opening <sup>50</sup> <sup>51</sup>. This is basically UI plumbing, but it's important for the flow: the user can spatially organize different threads or contexts of interaction. For example, one dot might bind to a Q&A panel with the AI, another dot to a brainstorming scratchpad panel. Clicking the dot brings that context up instantly. This spatialization means the user's attention can hop between contexts without losing them – effectively supporting *parallel attention streams*. Human attention in creative tasks often jumps around (think: jotting ideas, then verifying facts, then returning to writing). The Membrane UI's ability to keep multiple panels open or pinned (with Circle and Desk modes for minimizing) <sup>52</sup> <sup>53</sup> means the AI can be involved in all those micro-tasks concurrently. One panel might continuously summarize or rephrase text (like a live editor AI), another might generate new ideas on the side, all while the main writing continues. This concurrent assistance requires a well-managed interface to not overwhelm the user, hence the careful design of modes and the emphasis that even minimized panels remain easily accessible (large radial chips, etc.) <sup>54</sup>.

**Human-in-the-Loop Adjustments:** The interface also empowers the human to intervene in the generation process. For instance, if the AI starts going off-topic in a long answer, the user could stop generation and highlight a portion of the prompt or output and ask the AI a follow-up about it. The design encourages treating the AI as a collaborator – e.g., through "*self-prediction mirror*" snapshots <sup>55</sup> or "*intention prompts + attractors*" that the user sets and the AI respects <sup>56</sup>. An intention prompt might be something the user writes (like a goal description) and pins to the interface. The AI, being integrated with the UI, could always pay attention to that pinned intention (almost like a secondary prompt that stays constant). This way, even as the user and AI explore tangents, the system maintains a sense of the user's ultimate goal. In our context, that intention could be as simple as "write in a polite tone" or "aim for an audience of scientists". By having it in the interface (perhaps as a note on the Donuscope ring <sup>57</sup>), it influences all outputs without the user needing to repeat it every time.

In summary, the interface dynamics of the Donut/Membrane system show a concerted effort to make the boundary between human and AI **as coherent and rich as possible**. It is not just a text box for input and output; it's a multi-layered membrane with visual, temporal, and interactive affordances. The interface uses **geometric metaphors (toroids, spirals, radial patterns)** so that the user's experience remains anchored in the same conceptual universe as the theoretical model. This likely has cognitive benefits: the user internalizes the torus metaphor and can reason about their interaction in those terms (e.g., "maybe I need to add a bit of chaos to get out of this local minimum – let me turn up the dial or introduce an ambiguous prompt element"). The AI, on the backend, is just following algorithmic rules, but the interface translates some of those into perceptible signals for the human (like the spinning hologram for activity). Meanwhile, the human's nuanced decisions (like adjusting phrasing or toggling a UI option) are translated via the

membrane into signals the AI can use (perhaps altering system parameters or adding hidden tokens to the prompt context).

All these create a **feedback loop** that is tighter and richer than a plain exchange. We have the *outer loop* (human writes prompt, AI gives answer, human reads and writes a new prompt) and various *inner loops* (UI visual feedback guiding the human's next action, AI adjusting outputs based on persistent user-set intentions, etc.). When operating optimally, this feels like an almost fluid conversation where the human and AI are “in the zone” together. The Donut of Attention metaphor extends to this state as well – it’s when the two attentional fields (human and machine) synchronize, creating a larger joint field. The goal of the interface is to facilitate this synchrony, much as a good musical instrument disappears in the hands of a musician, allowing direct expression. In the ideal case, the Membrane UI becomes transparent: the human’s attention flows into the AI and back out, forming a continuous torus of creativity.

## Discussion

Our exploration has painted the picture of human prompt design and AI text generation as two interlocking components of a single cognitive system. By applying the Donut of Attention framework, we have highlighted deep parallels between human attention strategies and the AI’s attention mechanisms, as well as the importance of the interface that connects them. In this section, we reflect on the implications of these findings, address how ambiguity and creativity are handled in the system, and consider the broader significance of a *toroidal human-AI interface*.

One key insight is that **human-AI collaboration can be viewed as a coupled attentional field**. The human brings a lifetime of cognitive context, goals, and creativity, focusing it into a prompt. The AI brings a vast learned model of language, effectively a distilled collective intelligence, and focuses it via transformer attention on the prompt. When the interaction is successful, there is a sense in which the two attentions align – the AI “gets” what the human is asking and the human “gets” the AI’s output. The Donut metaphor helps conceptualize this as alignment of two toroidal fields into one larger torus (perhaps like two smoke rings merging). Achieving this alignment requires maintaining **coherence** in the interaction. We saw how coherence is maintained by designing prompts that are clear yet open, and by designing interfaces that provide guidance without imposition. Coherence does not mean rigidity; in fact, it includes the ability to handle divergence and novelty (the system can tolerate small contradictions or deviations as long as they eventually serve the overall goal, akin to the paraconsistent tolerance <sup>6</sup> ).

The handling of **ambiguity** in the interface is particularly noteworthy. Traditional software might try to eliminate ambiguity (every command must be exact). But in prompting, as we’ve discussed, ambiguity can be a feature, not a bug, when used deliberately. The AI model inherently represents uncertainty as probability distributions. Rather than forcing a single interpretation, the best prompts often let the model navigate multiple possibilities until it must choose (for example, a prompt might pose a mystery that the AI then resolves creatively). The Donut framework’s stance that we should “*embrace contradictions/ambiguity*” <sup>23</sup> is validated in practice by the richness of AI outputs to open-ended prompts. The interface could even amplify this by, say, offering multiple continuations for an ambiguous prompt and letting the user choose (a bit like giving the wavefunction’s several collapsed outcomes and having the human pick the desired branch). This would be an interesting extension: a *branching editor* where a single prompt spawns several parallel outputs in panels, and the user can pursue one or merge them. In effect, it would externalize the model’s internal branching (which normally happens implicitly via sampling) into a visible form. That would turn ambiguity from a one-time randomness into an interactive exploration tool.

Another aspect is how **creative exploration vs. focused execution** is balanced. The Creative Time mechanism and the concept of a *coherence corridor*<sup>46</sup> hint that there is an optimal zone between chaos and order for productivity. In prompting terms, this might correspond to using moderately high temperature or to alternating between divergent prompts (brainstorm something wild) and convergent prompts (organize these ideas). The interface's job is to help the user not get stuck in either extreme – neither in a chaotic flood of irrelevant text nor in a rigid question-answer that yields no surprises. The donut's edge-of-chaos is the sweet spot for creativity<sup>28</sup>. It is encouraging to see that system design explicitly considers this (e.g., the *Phase dials* in the design can bias suggestions toward more chaotic or more orderly<sup>46</sup>). Empirically, users find that if an AI's output is too deterministic it can become dull or repetitive, whereas if it's too random it loses coherence. Thus, a dynamic adjustment guided by either the system or the user leads to better outcomes. This is precisely the kind of higher-level control a well-designed interface can provide, abstracting away the fiddly parameters into intuitive controls (like a dial or mode toggle labeled perhaps "explore" vs "refine").

Our analysis also underscores the importance of **interpretable metaphors and visualizations** in human-AI systems. By using the same conceptual language (toroids, fractals, etc.) in the UI that we use to think about the cognitive theory, the system potentially shortens the learning curve for the user. They do not have to think in matrix multiplications or probabilities; they can think in terms of focusing, looping back, expanding out – terms that have visual correlates in the interface. For example, if a user feels "I need to refocus the AI on the main point," the interface might offer a literal action like clicking on the central bindu or highlighting the main point in the text and pressing a "focus" button. If they feel "let's explore a wild tangent," maybe toggling a "chaos monkey" spiral increases novelty. In a sense, the interface functions like training wheels for interacting with a complex model, but also as an extension of the user's mind. It externalizes what would otherwise be internal mental operations (like remembering to stick to a goal, or randomly brainstorming) into shared operations with the AI.

From the development perspective, the **Dev Journal** entries we cited provide a reality check for these high-level ideas. They show that implementing such an interface is non-trivial: issues like panel parenting, sync of rotations (for the donut visualization), and state persistence had to be solved<sup>58 59</sup>. It's clear that building a Membrane UI involved iterative fixes and close attention to user experience details (e.g., ensuring a search input is always visible with a halo so the user knows they can type<sup>60</sup>). The journal's collaborative notes where the developer uses an AI assistant ("web agent guidance") also illustrate our topic in action: the developer wrote prompts to an AI to get code suggestions, and commented on prompt clarity<sup>22</sup>. This meta-use of prompting to build a prompting interface is an example of the fractal nature of the whole endeavor – prompts to improve prompt-based systems! It highlights that as AI becomes part of development, the interface between human intention and code (which used to be programming languages) is also shifting to natural language. Ensuring that cognitive theories (like Donut of Attention) inform these tools could lead to more intuitive dev workflows as well, as seen with the "cooperation effectiveness" notes<sup>16</sup> aimed at continuous improvement of prompt techniques during development.

**Limitations and Future Work:** While the Donut of Attention provides a compelling holistic framework, it remains speculative in parts. Not all aspects of human cognition or AI processing can be neatly mapped to a torus or fractal. The metaphors should not be taken too literally; they are guides for design and understanding. For instance, one should be cautious in assuming a direct physical reality to these concepts (the attention field in a transformer is not actually a electromagnetic torus, it's a mathematical pattern – the metaphor is useful, not exact). Future empirical studies could test some implications: e.g., does using a visual toroidal interface measurably improve users' ability to maintain context in long conversations? Does

representing time in a circular form (as CTI does) improve adherence to work-break cycles and thereby output quality? These human-factor questions are ripe for exploration.

Another frontier is making the AI more aware of the metaphors so it can participate in the interface logic. For example, could the model maintain an internal “intention field” representation so that it doesn’t need to be reminded of the user’s high-level goals? Or could it detect when the interaction is drifting (losing coherence) and signal that on the interface (maybe the hologram flickers or a coherence meter drops)? Aligning the model’s own monitoring with the UI’s representation of state would effectively close the loop even tighter. Some modern research on AI alignment and self-monitoring (like transformers that can evaluate their own outputs for consistency) might integrate well with this.

Finally, the idea of *hybrid symbolic-geometric reasoning* <sup>19</sup> – which was mentioned in the theory – is exemplified in this interface. The human issues symbolic prompts, the AI processes geometrically in vector space, and the UI provides geometric visuals back to the human. It’s a cycle of translation between symbolic and geometric representations. This could be generalized beyond text generation: consider an AI that helps design actual geometry (like CAD software with an AI copilot). A similar interface approach might overlay symbolic instructions with geometric previews. Or in data analysis, a user could ask questions (symbolic) and get interactive charts (geometric) that they can manipulate, with the system maintaining an ongoing loop of understanding. The donut (or any cohesive metaphor) can serve as a common *language* between human intuition and machine computation across many domains.

## Conclusion

The process of going from a human’s prompt to an AI’s generated text is far more than a simple input-output transaction – it is an interplay of two systems of attention and cognition, bridged by a carefully crafted interface. By examining this process through the Donut of Attention framework, we have revealed how **human prompt design** is an act of mapping inner intent to language in a way that a machine can amplify, and how **AI text generation** is an act of focusing a learned distribution of language in response to that prompt. The toroidal, fractal, and holographic metaphors provide a unifying vision: the prompt and the output are connected by patterns that repeat across scales and by a shared “field” of context that both human and AI contribute to.

We have seen that effective prompts encapsulate a lot of meaning in a small space – like a **code point that expands into concentric shells** of content <sup>61</sup> – and that transformer models are designed to take such points and *blow them up* into detailed continuations, thanks to multi-headed attention distributing the prompt information throughout the network. We highlighted the importance of maintaining coherence without stifling creativity, using ambiguity as a feature and balancing on the edge of chaos for maximum creative potential. The system we analyzed (via its design docs and dev logs) explicitly embraces these ideas: from **paraconsistent UI design** that holds multiple ideas at once <sup>13</sup>, to **Creative Time** mechanisms that respect the natural rhythms of thought <sup>10</sup> <sup>8</sup>, to **visual metaphors** that keep the user oriented in a complex space of possibilities <sup>9</sup>.

In building the **cognitive interface** between human and AI, the Membrane UI demonstrates how thoughtful integration of theory and design can yield a workspace where the human and AI not only cooperate but seemingly start to **cothink**. When the interface is intuitive, the AI’s capabilities feel like a direct extension of the user’s mind – a extra set of cognitive hands, so to speak, operating in the background but under the guidance of the user’s intention. Conversely, the human becomes a crucial part

of the AI's functioning, providing direction, ethical grounding (deciding which outputs are acceptable), and novel information that the AI couldn't have on its own (like subjective preferences or real-time situational awareness).

Our journey through this prompt-to-text loop affirms that the *design of the prompt* and the *design of the interface* are as important to outcomes as the design of the AI model itself. Even a very advanced model benefits from a well-chosen prompt and a user who can effectively steer it – and both of those are fostered by a good interface that externalizes the right information and controls. The speculative yet scientific ethos of the Donut framework pushes us to imagine interfaces that are almost like **collaborative spaces**, where human and AI attention literally draw shapes and paths together (some researchers speak of “attention maps” shared with users). It might not be long before, for example, we have real neural network attention heatmaps shown in real-time as the model reads the prompt, or interactive knobs for different attention heads.

In conclusion, investigating the cognitive interface between human prompt design and AI generation reveals a rich, multi-layered collaboration. By framing it in terms of attention – a concept we can reason about both in cognitive and computational terms – we gain insights that translate into practical design principles. The toroidal Donut metaphor, as whimsical as it might sound, proves powerful in synchronizing ideas about focus, context, continuity, and creativity. It reminds us that great technology often works *with* human nature rather than against it: looping our attention in, not attempting to replace it. As AI systems become ever more prevalent, keeping the interface **human-centered** and **cognitively aligned** will be crucial. The prompt is where the human touch enters the system, and the more artfully we craft that touch – and design the surfaces it interacts with – the more the resulting human-AI symphony will achieve coherence, creativity, and perhaps even a touch of collective wisdom.

## References

1. **THEORETICAL\_FOUNDATIONS.md** – Internal conceptual notes defining the Donut of Attention framework (toroidal attention model, fractal-holographic logic, paraconsistency, etc.) 1 4 .
2. **UI\_ARCHITECTURE.md** – Design document for the Membrane UI architecture (shell modes, Circle/Desk radial menus, Creative Time integration, Solar Hologram overlay) 38 8 .
3. **dev\_journal.md** – Developer journal logs from the Donut project, detailing implementation decisions and collaborative notes (prompt clarity tracking, interface behavior fixes, integration of theoretical concepts) 16 22 .
4. **Palomares, I.** (2025). *How Transformers Think: The Information Flow That Makes Language Models Work*. KDnuggets. – Explains in accessible terms how transformer models tokenize input, use multi-headed attention to incorporate context, and generate text token-by-token 3 17 .
5. **Additional Project Notes** – Assorted imported summaries from retired design documents (e.g., Creative Time Maps, Holographic Attention notes) providing supporting context on UI features and theoretical alignments 10 28 .

3 5 17 24 25 26 27 29 30 How Transformers Think: The Information Flow That Makes Language Models Work - KDnuggets

<https://www.kdnuggets.com/how-transformers-think-the-information-flow-that-makes-language-models-work>

8 9 38 39 40 41 43 48 52 53 54 UI\_ARCHITECTURE.md

file:///file\_0000000066587246bf49dd4ae07982d4

16 22 37 44 45 50 51 58 59 60 dev\_journal.md

file:///file\_0000000051d071f4ae62b3c008adc811