

Cross-Modal Cognitive Mapping:

A New Framework for Understanding Human Thought Across Language and Vision

Author:

Ernan Hughes

Acknowledgment:

This work was developed with the assistance of AI cognitive tools as part of the methodology itself, reflecting the core principle of cross-modal human-AI collaboration.

Abstract

We propose a new framework for cognitive mapping that extends beyond traditional text embeddings by incorporating visual semantic generation and human selection behavior.

This method captures not only what users say, but how they **perceive and visualize** concepts — enabling deeper measurement of cognitive similarity and divergence across individuals, languages, and cultures.

Through a pipeline of text prompts → image generation → human selection → cross-user comparison, we construct a **cross-modal conceptual map** that reveals hidden dimensions of human thought inaccessible through text alone.

This framework has significant implications for universal concept understanding, cognitive anthropology, AI alignment, and next-generation multimodal intelligence systems.

1. Introduction

Modern AI systems primarily model user cognition through textual embeddings, mapping language into high-dimensional vector spaces.

However, human cognition is fundamentally multimodal: concepts are experienced not only through language but through **mental imagery, emotional tonality, texture, and subjective perception**.

Text embeddings alone cannot capture this richness.

We propose extending cognitive modeling into a **cross-modal framework**, combining:

- Natural language prompts,
- Visual semantic generation (images),
- Human selection behavior (preferred visualizations),
- Cross-user cognitive resonance and divergence mapping.

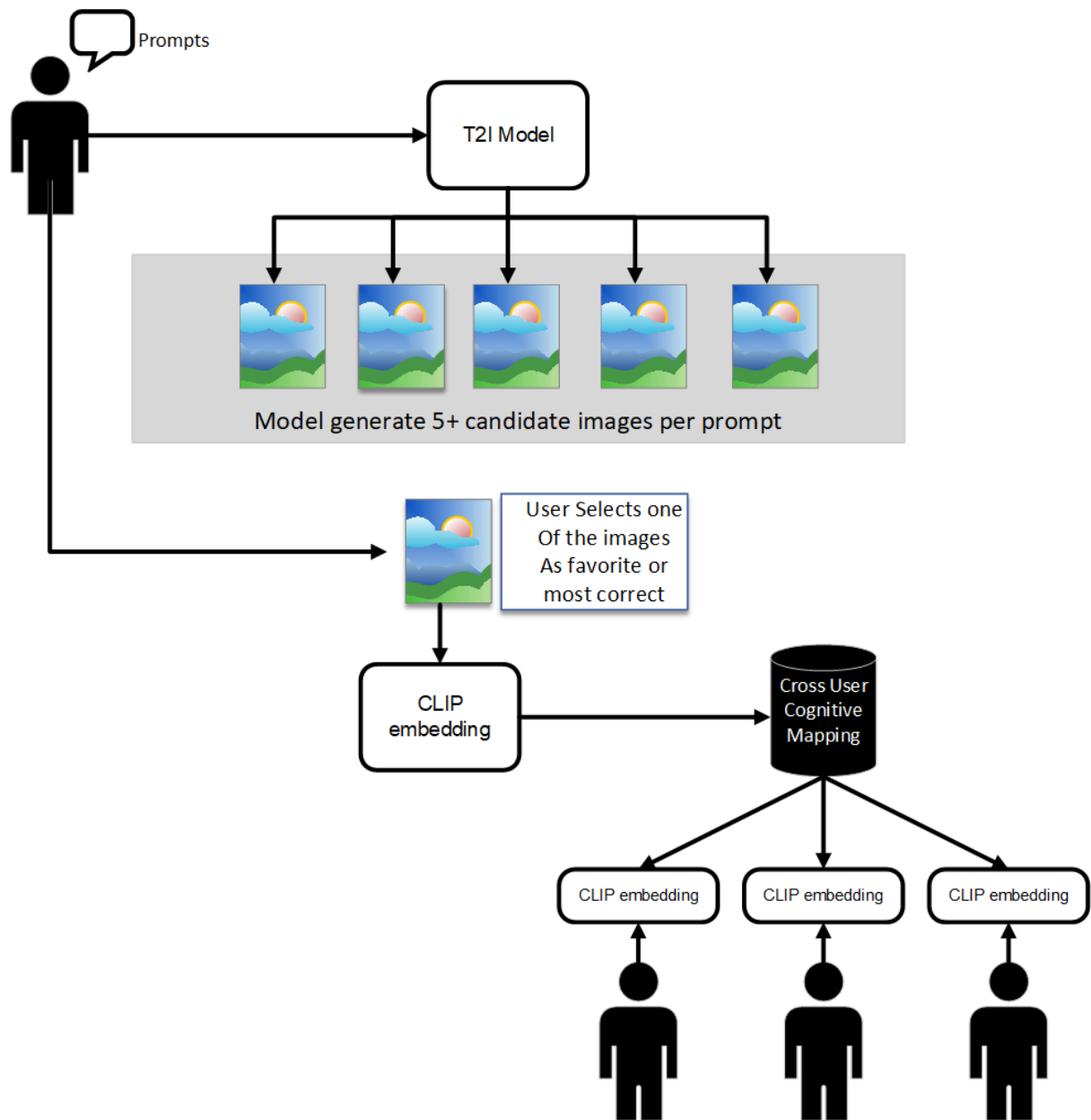
By observing how users imagine and select among visualizations of their own thoughts, we can map their cognitive architecture far beyond what text models alone reveal.

2. Method

The proposed method operates across four layers:

Layer	Description
Text Space	Capture initial user prompts in natural language.
Image Space	Generate multiple visual candidates per prompt using text-to-image (T2I) models.
Selection Space	Let humans select preferred visualizations, capturing subjective cognitive styles.
Cognitive Comparison Space	Embed selected images and compare across users to map convergence and divergence of conceptualization.

Process Pipeline:



3. Applications

The cross-modal cognitive mapping framework enables new forms of analysis and interaction:

Domain	Opportunity
Cross-Language Cognitive Bridging	Identify shared concepts across users with different native languages through visual semantic similarity.
Cognitive Anthropology at Scale	Observe and analyze conceptual patterns across cultures, regions, and generations.
Personalized AI Alignment	Train AI agents to align not only with user language, but with their internal imagery and conceptual styles.
Concept Evolution Tracking	Study how human perceptions of abstract concepts evolve over time (e.g., "home", "freedom", "future").
Memory Systems and Visual Diaries	Build tools for individuals to store and explore their internal thought journeys visually, not just verbally.

4. Case Example: Divergent Visualization

Prompt: "Describe a peaceful home."

User	Top Selected Image	Latent Visualization
User A	Wooden cabin by snowy mountains	Solitude, cold serenity
User B	Beach hut with palm trees	Warmth, openness

Despite using nearly identical language, the users' internal conceptions of "peace" diverge significantly — revealing hidden cognitive differences invisible to text-only models.

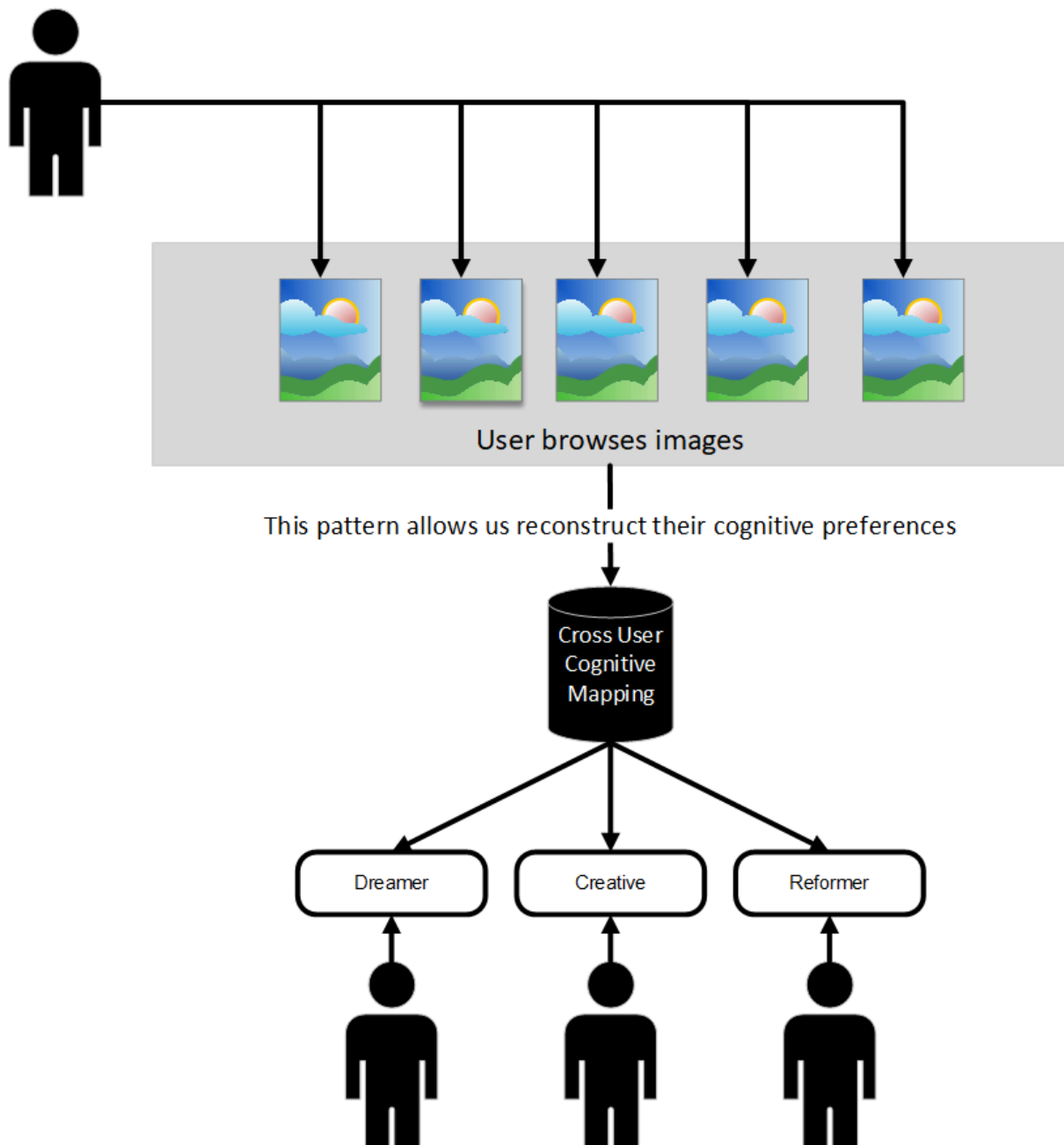


5. Reverse Cognitive Mapping

Beyond forward mapping (prompt → image → selection → map), this framework enables **reverse inference**:

By observing a user's path through a landscape of images (clicks, hovers, skips, linger time), and embedding the images they engage with, we can **reconstruct their latent cognitive preferences** without requiring explicit prompts.

Thus, the system supports **both active** (prompt-driven) **and passive** (behavior-driven) **cognitive modeling**.



6. Cognitive Atlas (Future Work)

To structure reverse cognitive inference, a **Cognitive Atlas** must be constructed:

- Define cognitive archetypes (e.g., dreamers, builders, explorers, philosophers).
- Map their characteristic text-image-selection patterns.
- Train models to match new users against these reference trajectories.

This will enable faster and more accurate cognitive type prediction from behavioral traces alone.

7. Conclusion

By extending cognitive modeling across text, image generation, and human selection, we unlock deeper access to the conceptual architectures of the human mind.

This cross-modal cognitive mapping approach moves AI understanding beyond language, toward **perception**, **imagination**, and **shared internal worlds**.

We believe this method opens new frontiers in AI-human collaboration, cultural anthropology, personalized alignment, and cognitive self-discovery.
