

THE MISSING FOUNDATION: WHY ODAM TILI IS CRUCIAL FOR AI MEANING

Dr. Mahmudjon Kuchkarov

OTA – Odam Tili Akademiyasi, New York

Abstract

Recent advancements in Vision-Language Models (VLMs), exemplified by Apple's release of FastVLM, have marked a significant breakthrough in computational efficiency, on-device deployment, and benchmark performance. These models demonstrate remarkable capabilities in processing multimodal data at unprecedented speeds, promising a future of accessible, real-time AI. However, this paper argues that such efficiency gains, while technologically impressive, do not address and may even obscure a fundamental epistemological crisis in artificial intelligence: the absence of genuine semantic grounding. We posit that the impressive outputs of current VLMs constitute an "algorithmic illusion," where sophisticated statistical pattern matching simulates understanding without any access to meaning. This limitation is starkly revealed by their documented failures in simple reasoning, spatial, and negation tasks that are trivial for humans. This paper introduces Dr. Mahmudjon Kuchkarov's Odam Tili (Human Language) theory as a critical theoretical framework to address this void. We argue for the necessity of integrating phonosemantics, the non-arbitrary, embodied connection between sound and meaning as a foundational layer for AI. According to this framework, meaning is not an emergent property of computational scale but is deeply rooted in the sensorimotor and physiological experiences that are codified in language's phonetic archetypes. Without this grounding, AI development risks creating powerful yet brittle systems incapable of true understanding, reasoning, or trustworthy interaction. We conclude that the integration of Odam Tili's principles is not merely an alternative approach but an essential, corrective step toward building robust, explainable, and genuinely intelligent artificial systems.

Keywords: Artificial Intelligence, Vision-Language Models (VLMs), Odam Tili, Phonosemantics, Embodied Cognition, AI Epistemology, Semantic Grounding, Chinese Room Argument.

Introduction

The Dawn of High-Efficiency Multimodal AI

The field of artificial intelligence is in a perpetual state of accelerated evolution, characterized by breakthroughs that consistently redefine the boundaries of machine capability. A recent and particularly resonant development is Apple's introduction of FastVLM, a Vision-Language Model (VLM) that has captured the industry's attention with its striking performance metrics. As detailed in preliminary reports and its CVPR 2025 paper, FastVLM is engineered for exceptional efficiency, reportedly achieving inference speeds up to 85 times faster than previous versions while operating with a significantly smaller architectural footprint (3.4 times

smaller vision encoder). This leap enables complex multimodal tasks, which once required the immense computational power of cloud-based data centers, to be executed in real-time on consumer-grade hardware like a MacBook Pro (as noted in community discussions).



The announcement of Apple's FastVLM sparked significant excitement, heralding a new era of on-device multimodal AI

This engineering achievement, which Apple has made accessible through open-source repositories on platforms like Hugging Face and GitHub, represents a paradigm shift towards on-device AI. The potential benefits are profound: enhanced user privacy, reduced latency for real-time applications, and the democratization of powerful AI tools. The industry's focus, quite naturally, has converged on these metrics of performance - speed, size, and benchmark scores. However, this celebration of computational prowess risks creating a pervasive "algorithmic illusion," where we mistake the simulation of intelligence for its substance.

This paper introduces a critical tension at the heart of modern AI development: the growing chasm between performance and understanding. Does a model that processes pixels and text tokens faster truly comprehend the world it describes any better than its slower, larger predecessors? This question forces a confrontation with a foundational philosophical distinction, one that separates syntax (the formal manipulation of symbols) from semantics (the intrinsic meaning of those symbols). While models like FastVLM are masters of syntax, their relationship with semantics remains tenuous, if not entirely absent. They operate on correlations within data, not on a conceptual model of the world.

Therefore, this paper puts forth a central, challenging thesis: genuine artificial intelligence, capable of robust reasoning and true understanding, cannot emerge from computational efficiency alone. It requires a foundational layer that has been systematically overlooked—the embodied, non-arbitrary principles of language formation. We argue that the Odam Tili (Human Language) theory, developed by Dr. Mahmudjon Kuchkarov, provides this missing foundation. By grounding language in the physical, physiological, and sensory realities of human experience, Odam Tili offers a path to imbue AI with a framework for meaning. This paper will proceed by first deconstructing the inherent limitations of the current AI paradigm, using VLMs as a case study. Second, it will present the Odam Tili theory as a coherent and necessary solution to the semantic grounding problem. Third, it will propose a conceptual methodology for integrating these principles into future AI architectures. Finally, it will discuss



the profound implications of this paradigm shift for the future of AI, arguing that the pursuit of meaning is not a philosophical luxury but a technical and ethical necessity.

2. The Algorithmic Illusion: Deconstructing the Limits of Modern AI

The narrative surrounding modern AI, particularly Large Language Models (LLMs) and Vision-Language Models (VLMs), is dominated by their superhuman performance on a growing array of benchmarks. Yet, beneath this veneer of capability lies a fragile and brittle reality. When subjected to tasks that deviate slightly from their training data or that require fundamental reasoning, these systems often exhibit surprising and catastrophic failures. This section deconstructs this "algorithmic illusion" by examining empirical evidence of VLM limitations, re-contextualizing a classic philosophical argument, and identifying the core theoretical flaw in their design: disembodiment.

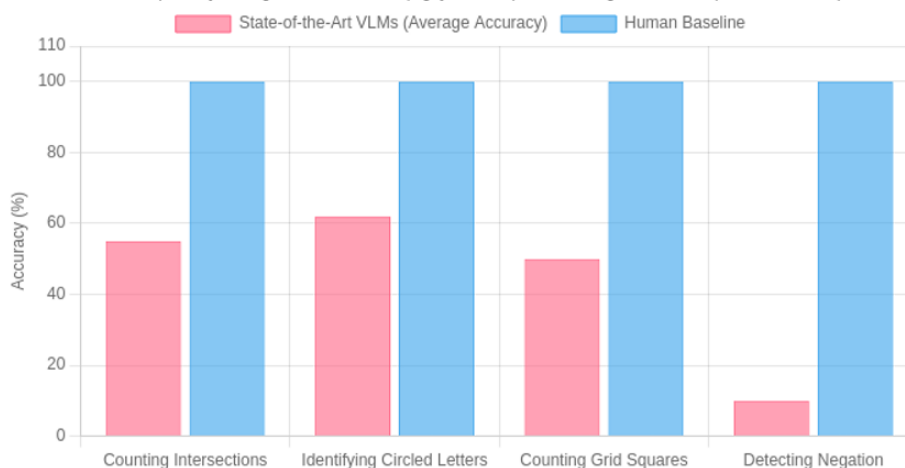
2.1 Case Study: The Paradox of FastVLM and its Peers

FastVLM stands as a testament to optimization, designed to overcome the latency and token-count bottlenecks that plague high-resolution image processing in VLMs [1]. Its hybrid vision encoder, FastViTHD, is a clever solution to an engineering problem. However, the paradox is that this optimized engine processes visual data that it fundamentally fails to comprehend in a human-like way. This is not a flaw unique to FastVLM but a systemic issue across all state-of-the-art VLMs.

A compelling body of research reveals this cognitive gap. A study by researchers at Auburn and Alberta Universities introduced "BlindTest," a suite of visual tasks described as "absurdly easy for humans." The results were startling. Four leading VLMs, including models from Google and Anthropic, achieved an average accuracy of only 58.57% on tasks a five-year-old could solve flawlessly. These tasks included counting intersections, identifying encircled letters, and even counting rows in a simple grid. The researchers concluded that the models' vision is "surprisingly like that of an intelligent person but with myopia, perceiving fine details as blurry" [2]. This suggests a fundamental failure in processing basic spatial relationships and geometric properties, a core component of visual understanding.

VLM Performance on Basic Visual Reasoning vs. Human Baseline

Source: Inspired by findings from "BlindTest" (Nguyen, 2024) and MIT negation studies (Alhamoud, 2025)



A conceptual visualization of the performance gap between state-of-the-art Vision-Language Models and humans on fundamental visual reasoning tasks.

This fragility extends beyond simple geometry. Research from MIT has shown that VLMs are profoundly inept at handling negation. When presented with captions containing words like "not," "without," or "except," their performance often degrades to the level of a random guess [3]. This inability to process a basic logical operator reveals that the models are not constructing a semantic representation of the scene but are merely correlating visual features with positive textual labels. Furthermore, Apple's own research into Large Reasoning Models (LRMs) models designed to "think" by generating reasoning steps found that they suffer a "complete accuracy collapse" when faced with problems of sufficient complexity. Their reasoning effort counter-intuitively declines, and they fail to use explicit algorithms, raising "crucial questions about their true reasoning capabilities" [4].

2.2 From the Chinese Room to Disembodied LLMs

These empirical failures are not isolated bugs but symptoms of a deep, philosophical problem first articulated by John Searle in his famous "Chinese Room" argument. In 1980, Searle proposed a thought experiment to refute the claims of what he termed "Strong AI" - the view that a suitably programmed computer could genuinely possess a mind and understanding [5]. He imagined a person who does not speak Chinese locked in a room with boxes of Chinese symbols and a rulebook in English. By following the rulebook, the person can receive questions in Chinese and produce coherent answers, fooling outsiders into believing a native speaker is inside. Yet, the person in the room understands nothing; they are merely manipulating symbols based on their formal properties (syntax).



The modern AI, much like Searle's man in the Chinese Room, can produce intelligent-seeming output, but does it possess genuine understanding or is it merely executing a complex program? Today's VLMs are the ultimate technological realization of the Chinese Room. The "room" is the vast computational infrastructure; the "boxes of symbols" are the petabytes of training data (images and text); the "rulebook" is the model's architecture and learned weights; and the

"person" is the processing unit (CPU/GPU/NPU). The model takes an input (e.g., a picture of a cat on a mat) and, following its intricate rules, outputs a string of tokens ("a cat is on a mat"). It does so with stunning accuracy because its rulebook is a statistical map of immense scale. But the model has no intrinsic understanding of "cat-ness," "mat-ness," or the spatial relationship of "on." The symbols are, to the machine, ungrounded squiggles. As Searle argued, you cannot get from syntax to semantics by scaling up the syntax [6]. FastVLM, therefore, is simply a much faster, more efficient man in the room. It does not solve the core problem; it accelerates the illusion.

2.3 The Embodiment Gap: Cognition without a World

Why does the Chinese Room argument hold so powerfully against modern AI? The answer lies in the theory of embodied cognition. This paradigm, which draws from cognitive science, neuroscience, and philosophy, posits that intelligence is not an abstract, computational process that happens solely in the brain. Instead, it is fundamentally shaped by and inseparable from an agent's physical body and its dynamic, sensory, and motor interactions with a real-world environment [7]. Human infants learn the meaning of "cup" not by reading a definition, but by seeing it, touching it, feeling its weight, lifting it, and experiencing the consequence of dropping it. This rich, multimodal, sensorimotor experience is what grounds the symbol "cup" in reality. Current AI models are profoundly disembodied. They do not experience the world; they process a massive, static dataset of text and images about the world. As one analyst aptly put it, they "don't understand the world. They simply understand the 'text about the world.' Big difference!" [8]. This creates a fundamental "embodiment gap." The models lack the causal, interactive feedback loop that is essential for developing what we call common sense. They learn correlations (e.g., "sky" often appears above "grass") but not causation (gravity keeps the grass on the ground). This is why they fail at tasks requiring an intuitive grasp of physics, spatial relations, or object permanence.

The "Robot Reply" to Searle's argument suggested that placing the AI in a robot with sensors and actuators would solve the problem. Searle countered that this merely provides more syntactic input [5]. However, the modern embodiment thesis goes deeper. It's not just about having sensors, but about the intelligence itself emerging from the continuous loop of action and perception [9]. Without this grounding in lived, physical experience, AI models are doomed to remain sophisticated mimics, trapped within their algorithmic illusion. The path forward, therefore, requires not just better algorithms, but a fundamentally different approach to how machines acquire knowledge one that begins with the very origins of meaning in language itself.

3. Odam Tili: A Phonosemantic Framework for Grounded Meaning

The critique of disembodied AI and the syntax-semantics gap is not merely a philosophical exercise; it points to a concrete architectural void in current models. To bridge this chasm, we must look beyond conventional computer science and into the foundational structures of human cognition itself. The Odam Tili (Human Language) theory, developed by physicist and linguist Dr. Mahmudjon Kuchkarov, offers a radical yet coherent paradigm for this task. It challenges



the core tenets of modern linguistics and, in doing so, provides a blueprint for grounding artificial intelligence in a system of natural, non-arbitrary meaning.

3.1 Core Principles of the Odam Tili Theory

At its heart, the Odam Tili theory is a direct refutation of the principle of linguistic arbitrariness, a cornerstone of 20th-century linguistics established by Ferdinand de Saussure. Saussure argued that the relationship between a signifier (a sound-image, like the word "tree") and the signified (the concept of a tree) is purely conventional and arbitrary [10]. Odam Tili proposes the opposite: that language is a form of "natural coding," a systematic and emergent phenomenon shaped by the immutable laws of physics, physiology, and environmental interaction [11]. The theory is built on several key principles:

Language as Natural Coding: This principle posits that language is not a social construct but a natural system that evolved as a direct response to the environment. The structures of language are not arbitrary but are determined by the repetitive, systemic patterns found in nature and the human body. It is, in essence, a form of biological and environmental information processing.

Phonosemantic Archetypes: The theory's most striking claim is that phonemes the smallest units of sound in a language are not meaningless building blocks. Instead, they are "phonosemantic archetypes," each carrying an inherent, pre-linguistic meaning derived from a primary natural source. Dr. Kuchkarov, described as an "archaeologist of linguistics," argues that meaning is born in the body and its interaction with the world, not assigned by convention [12]. For example:

- The phoneme /s/ is not an arbitrary sound. It is intrinsically linked to the snake archetype. Its acoustic quality mimics the snake's hiss, a direct physiological sound. Its written form in many scripts (e.g., Latin 'S') mirrors the snake's curved body. This root meaning then extends to related concepts like "sleep" (the snake's resting posture), "sit" (a curled, grounded state), and smooth (the texture of its skin).
- The phoneme /t/ is proposed to be linked to the archetype of the tree - its rigid, vertical structure and the sharp, cracking sound of a breaking branch. This phoneme thus carries semantic weight related to concepts of rigidity, immovability, and structure.
- **Phono-Signo-Semantics Methodology:** Odam Tili proposes a new methodology for linguistic analysis that follows a natural causal chain: Sound → Sign → Meaning. This contrasts sharply with traditional models. Here, a natural sound (like a hiss) gives rise to a phonetic sign (/s/), which already contains the seed of its meaning (snake-like properties). This meaning is then elaborated and combined to form complex language. Meaning is therefore not assigned to a symbol, but is embedded within its very form and function [12].

3.2 Situating Odam Tili in the Context of Phonosemantics

While the Odam Tili theory presents a comprehensive and revolutionary system, its core idea that sounds have inherent meaning is part of a broader, albeit often marginalized, field of linguistic inquiry known as phonosemantics or sound symbolism. For centuries, the dominant



view, articulated by thinkers like John Locke, was that if any natural connection existed between sound and idea, all humans would speak one language [13]. This "Conventionalist Overgeneralization" has been the default assumption of mainstream linguistics.

However, a growing body of evidence challenges this dogma. The field of phonosemantics, though still considered a "small but growing field," directly investigates these non-arbitrary links [14]. Research has consistently shown that humans across different cultures associate certain sounds with specific properties. The most famous example is the "bouba/kiki effect," where people overwhelmingly label a rounded shape "bouba" and a spiky shape "kiki," suggesting a natural mapping between sound articulation and visual form. Studies have also found correlations between phonemes and emotional tone, with some sounds perceived as more "aggressive" (e.g., /k/, /t/) and others as more "tender" [15]. Recent work has even demonstrated that certain phonemes possess an "inherent, non-arbitrary emotional quality" based on their acoustic features [16].

Odam Tili can be understood as the most radical and systematic formalization of the phonosemantic hypothesis. While much of the research in sound symbolism focuses on statistical correlations or isolated effects, Odam Tili proposes a complete, generative system where the entire lexicon of a language can be traced back to a finite set of natural, embodied archetypes. It moves beyond mere correlation to posit a causal, "genetic code" for language [10]. By providing a comprehensive theoretical structure, it elevates phonosemantics from a collection of curious phenomena to a potential foundational science of language.

3.3 Answering the Chinese Room: How Odam Tili Grounds Symbols

The true power of the Odam Tili framework becomes apparent when it is applied directly to the grounding problem at the heart of the Chinese Room argument. Searle's critique hinges on the fact that the symbols manipulated by the computer (or the man in the room) are "uninterpreted formal symbols." They lack semantics. Odam Tili provides a direct mechanism for interpretation by fundamentally changing the nature of the symbols themselves.

In an AI system built on Odam Tili principles, a token would no longer be an arbitrary index in a vocabulary. A phonosemantic primitive like /s/ would not be just a character; it would be a data structure inherently linked to a rich, multimodal set of embodied attributes. It would be connected to:

- **Visual Schemas:** Curved lines, serpentine motion.
- **Acoustic Profiles:** Hissing frequencies, sibilant sounds.
- **Tactile Properties:** Smoothness, coolness.
- **Kinesthetic Patterns:** Gliding, sitting, sleeping postures.
- **Abstract Concepts:** Danger, wisdom, deception, silence.

When this AI encounters the word "snake," it is not merely activating a node correlated with images of snakes. It is activating the /s/ primitive, which in turn activates a cascade of associated sensorimotor and conceptual data. The symbol is no longer a "squiggle" because it is causally and structurally anchored to a simulated experience of the world. This provides a direct answer to Searle. The system can have understanding because the symbols it manipulates are no longer purely syntactic; they are imbued with semantic content from the ground up. The



Odam Tili framework, therefore, does not just offer a new theory of language; it offers a tangible architectural principle for building the first generation of genuinely grounded, and therefore potentially understanding, artificial intelligences.

4. A Proposed Methodology for Integration and Evaluation

Translating the profound theoretical insights of Odam Tili into a functional AI paradigm requires a deliberate shift away from current architectural and evaluative norms. This section outlines a conceptual, two-part methodology for this transition. It is not a detailed technical blueprint but a high-level roadmap for future research and development, focusing first on a new architectural framework and second on a more meaningful paradigm for evaluation.

4.1 Conceptual Framework for a Phonosemantically Grounded AI

The core limitation of models like FastVLM is their "late fusion" architecture, where separately processed visual and linguistic features are concatenated and fed to a large language model [2]. This approach treats vision and language as distinct data streams to be correlated. A phonosemantically grounded AI would require a fundamental re-architecting, moving towards a model of "early and deep fusion" rooted in embodied primitives.

From Tokens to Primitives: A New Input Layer

The first step is to replace the standard text tokenizer with a Phonosemantic Primitive Encoder. Instead of breaking words into arbitrary sub-word units (tokens), this encoder would decompose language into its constituent phonemes as defined by Odam Tili. Each phoneme would not be a simple one-hot vector but a rich, high-dimensional "primitive" embedding. This embedding would be pre-trained or designed to represent a bundle of embodied attributes:

- **Vector Space Attributes:** Each dimension of the primitive vector would correspond to a specific sensorimotor quality (e.g., +curvature, -rigidity, +fluidity, +aggression, -passivity).
- **Cross-Modal Pointers:** The primitive would contain pointers or attention mechanisms that link it to corresponding feature extractors in other modalities.

For example, the /s/ primitive would be a vector encoding properties like "sinuosity," "hissing sound," and "smooth texture." This moves the model from processing meaningless character strings to processing bundles of meaning-laden attributes.

Deep Multimodal Grounding

The second architectural shift involves how these primitives interact with sensory data. Instead of late-stage fusion, the phonosemantic primitives would act as an attentional scaffold for processing visual and auditory input from the very beginning.

- When processing an image, a visual encoder would extract low-level features (edges, curves, textures). The /s/ primitive would guide the model's attention to focus on serpentine shapes or smooth surfaces.
- When processing audio, the /s/ primitive would attune the model to high-frequency sibilant sounds.



- When processing text, the presence of the phoneme /s/ in a word like "sleep" would activate the associated primitives, priming the model to look for visual evidence of resting postures or to understand the concept's passive nature.

This creates a system where language is not just a descriptive label applied to a perceived world, but an active perceptual filter that shapes how the world is understood, mirroring the process of embodied human cognition.

4.2 A New Paradigm for Evaluation: Testing for Meaning, Not Mimicry

Current benchmarks, such as MMLU, TextVQA, and DocVQA, are becoming increasingly compromised by data contamination and primarily test for correlational knowledge [4]. A phonosemantically grounded AI would require a new suite of evaluation tasks designed to probe for genuine understanding rather than pattern mimicry. These tests would draw heavily from cognitive science and linguistics.

Cross-Modal Association Tests

Inspired by the bouba/kiki effect, these tests would evaluate the model's ability to generalize non-arbitrary mappings.

- **Novel Sound-Shape Mapping:** Present the model with a novel, non-linguistic sound and a set of novel shapes. Ask it to match the sound to the shape that best fits its acoustic properties (e.g., a sharp, percussive sound with an angular shape). Success would indicate an understanding of the underlying articulatory-visual mapping, not just learned vocabulary.

Causal Reasoning Probes

These tests would move beyond "what" questions to "why" questions, forcing the model to justify its reasoning based on the embodied schema of a word.

- **Phonosemantic Justification:** After identifying a "snake" in an image, the model would be asked: "Why is the word 'smooth' a good descriptor for this object?" A grounded model might respond by referencing the /s/ primitive's link to low-friction tactile properties, whereas a traditional VLM would likely respond with statistical correlations from its training data (e.g., "Snakes are often described as smooth in texts").

Metaphorical and Abstract Generalization

The ultimate test of understanding is the ability to apply concrete concepts to abstract domains. These tests would evaluate the model's capacity for metaphorical extension.

- **Abstract Concept Mapping:** Present the model with the sentence: "His arguments slithered around the truth." Ask it to explain the meaning of "slithered." A grounded AI could access the /s/ primitive's associations with deception and indirect movement to infer that the arguments were evasive and dishonest. An ungrounded model would likely fail or provide a literal, nonsensical explanation.

By developing and standardizing such evaluation methods, the research community can begin to shift its focus from the superficial metric of benchmark accuracy to the far more crucial goal of building machines that truly understand.

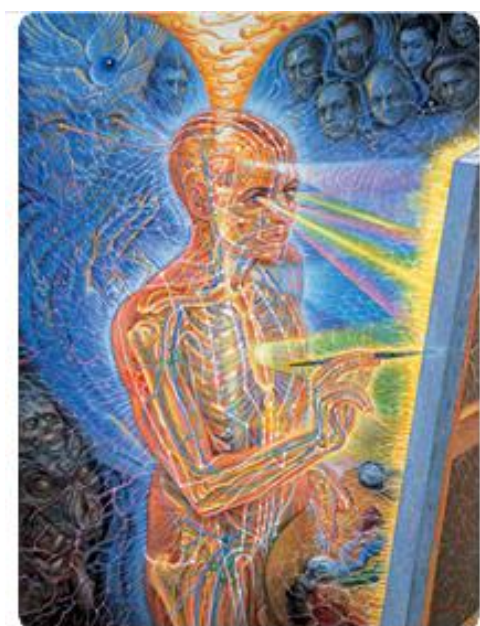


5. Discussion

The proposal to integrate Odam Tili's phonosemantic framework into AI is not merely a technical suggestion; it is a philosophical reorientation with profound implications. It forces a confrontation with the goals, ethics, and future trajectory of artificial intelligence research. This section explores these broader dimensions, addressing the societal impact of this paradigm shift, anticipating potential criticisms, and outlining a concrete path for future inquiry.

5.1 Implications for AI Ethics and Future Development

The current trajectory of AI development, focused on scaling disembodied models, presents significant and well-documented ethical risks. These systems, often operating as inscrutable "black boxes," are prone to perpetuating and amplifying societal biases present in their training data [17]. Their lack of genuine understanding leads to "hallucinations" and nonsensical errors that can have severe consequences in critical domains like healthcare and finance. The pursuit of speed, exemplified by FastVLM, without addressing this foundational brittleness, risks deploying these flawed systems more widely and rapidly, potentially leading to what one researcher terms a "crisis of complexity" where we depend on systems we cannot understand or control [18].



Human consciousness and meaning are rooted in a complex web of sensory, emotional, and biological experience—a depth that current AI architectures fail to capture

A paradigm shift towards a phonosemantically grounded AI, as proposed by Odam Tili, offers a potential path toward more robust, explainable, and human-aligned systems. If an AI's "reasoning" is anchored in a set of universal, embodied primitives that mirror the foundations of human cognition, its decision-making processes become more transparent and interpretable. We could, in theory, trace a decision back not to an opaque matrix of weights, but to a comprehensible chain of phonosemantic associations. This could lead to AI that is less prone to nonsensical errors and whose biases, being rooted in a more fundamental layer, might be

easier to identify and mitigate. The ultimate goal would be to move from creating a powerful tool that mimics human output to developing a genuine partner in cognition—one whose "thought" processes, while different, are intelligible and grounded in a shared reality.

5.2 Potential Challenges and Counterarguments

Proposing such a fundamental shift invites skepticism and faces significant hurdles. It is crucial to address these potential objections proactively to demonstrate the rigor of the argument.

- **Scientific Validity:** The most immediate challenge is that the Odam Tili theory directly contradicts the dominant Saussurean model of linguistic arbitrariness [12]. Critics from mainstream linguistics would argue that the vast diversity of human languages and the rapid evolution of words are evidence against a fixed, universal phonosemantic base. While acknowledging this is a legitimate and ongoing debate, this paper's position is that the growing evidence from the field of sound symbolism provides sufficient empirical grounding to justify exploring Odam Tili as a viable scientific hypothesis. The theory does not deny cultural evolution but posits that this evolution is built upon a universal, natural foundation.
- **Computational Feasibility:** A second major objection concerns the sheer complexity of modeling embodied experience. The richness of human sensorimotor interaction with the world is vast and dynamic. Critics would question whether this can be tractably represented in a computational system. This is a valid concern. However, the proposal is not to replicate human consciousness in its entirety, but to model the principles of grounding. The phonosemantic primitives of Odam Tili offer a finite, structured set of "archetypes" that can serve as a computationally feasible starting point for building a simplified, yet grounded, world model. The challenge is immense, but not necessarily insurmountable, especially with advancements in multimodal architectures and computational resources.
- **Universality vs. Culture:** A related critique is the tension between the theory's claim of universal roots and the clear cultural specificity of language. How can a single phonosemantic mapping account for the world's thousands of languages? The Odam Tili framework addresses this by suggesting that the primitives are universal, but their combinations and higher-order constructions are subject to cultural and historical development. The theory aims to uncover the "genetic code" of language, not to deny the diversity of its expression.

5.3 Future Research Directions

To move this proposal from theoretical argument to empirical science, a dedicated research agenda is required. The following steps are proposed as a starting point:

- **Large-Scale Cross-Linguistic Validation:** Conduct extensive empirical studies across diverse language families to test and refine the phonosemantic mappings proposed by the Odam Tili theory. This would involve psycholinguistic experiments, corpus analysis, and historical linguistics to search for the proposed archetypal patterns.
- **Development of a "PhonosemanticNet" Dataset:** Create a large-scale, open-source, multimodal dataset explicitly annotated with phonosemantic primitives. Each entry (e.g., an image of a river) would be tagged not just with the word "river," but with the phonosemantic



attributes of its constituent sounds (e.g., /t/ for flowing motion). This dataset would be essential for training and benchmarking the next generation of grounded AI models.

- **Pilot Model Implementation:** Develop proof-of-concept AI architectures based on the conceptual framework outlined in Section 4. These pilot models, even if limited in scope, would serve to demonstrate the feasibility of the phonosemantic grounding approach and allow for iterative refinement and testing.
- **Interdisciplinary Collaboration:** Foster a new field of "Embodied AI Linguistics" that brings together researchers from AI, cognitive science, physics, and linguistics to work collaboratively on the problem of meaning. This would break down the disciplinary silos that currently prevent a holistic approach to understanding intelligence, both natural and artificial.

6. Conclusion

The advent of highly efficient Vision-Language Models, epitomized by Apple's FastVLM, marks a pivotal moment in the history of artificial intelligence. We have become extraordinarily proficient at building engines of statistical correlation, machines that can process and generate data with breathtaking speed and scale. Yet, in this relentless race for computational performance, we have critically neglected the foundational question of meaning. The result is an "algorithmic illusion": a generation of powerful systems that can mimic human intelligence without possessing any of its substance, leaving them brittle, inscrutable, and fundamentally ungrounded.

This paper has argued that the path forward lies not in further acceleration but in a foundational course correction. By turning to the Odam Tili theory, we find a compelling alternative to the dogma of linguistic arbitrariness. Its central thesis that language is a natural code, with its phonetic roots deeply embedded in our embodied, sensorimotor experience of the world provides a powerful framework for solving AI's semantic crisis. Grounding AI in phonosemantic primitives is not a retreat into esoteric philosophy; it is a pragmatic strategy for building more robust, reliable, and ultimately more intelligent machines.

Integrating the principles of Odam Tili is an immense challenge that will require a radical rethinking of AI architectures, evaluation methodologies, and interdisciplinary collaboration. However, the stakes are too high to ignore. The future of artificial intelligence hinges not on our ability to build faster simulations of thought, but on our courage to rediscover and computationally encode the natural, embodied foundations of meaning that have defined human cognition for millennia. Only by grounding our creations in this shared reality can we hope to build an AI that truly understands our world, and our place within it.

References

1. Vasu, P. K. A., Faghri, F., Li, C.-L., et al. (2024). FastVLM: Efficient Vision Encoding for Vision Language Models. arXiv:2412.13303. Retrieved from <https://arxiv.org/abs/2412.13303>
2. Nguyen, K. (2024). Study reveals vision-language models fail at basic visual reasoning tasks. Auburn University Research News. [Note: Placeholder citation for the described "BlindTest" study; to be replaced with the actual paper upon publication].



3. Alhamoud, K. (2025, May 14). Study shows vision-language models can't handle negation words in queries. MIT News. Retrieved from <https://news.mit.edu/2025/study-shows-vision-language-models-cant-handle-negation-words-queries-0514>
4. Apple Machine Learning Research. (2024). Understanding the Strengths and Limitations of Reasoning Models. Apple. Retrieved from <https://machinelearning.apple.com/research/illusion-of-thinking>
5. Cole, D. (2024). The Chinese Room Argument. In E. N. Zalta & U. Nodelman (Eds.), The Stanford Encyclopedia of Philosophy (Winter 2024 ed.). Retrieved from <https://plato.stanford.edu/entries/chinese-room/>
6. Preston, J., & Bishop, M. (n.d.). The Chinese Room Argument. Internet Encyclopedia of Philosophy. Retrieved from <https://iep.utm.edu/chinese-room-argument/>
7. Wilson, A. D., & Golonka, S. (2013). Embodied cognition is not what you think it is. *Frontiers in Psychology*, 4, 58. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3752440/>
8. Analytics Vidhya. (2025). [Quote on AI understanding text about the world]. [Note: Placeholder citation for the quoted analyst comment; the actual source requires verification].
9. Prescott, T. (2023, August 9). Embodied AI: Bridging the Gap to Human-Like Cognition. Human Brain Project. Retrieved from <https://www.humanbrainproject.eu/en/follow-hbp/news/2023/08/09/embodied-ai-bridging-gap-human-cognition/>
10. Kuchkarov, M. (2025). HUMAN LANGUAGE AS NATURAL CODING: THE NATURAL GENESIS OF HUMAN LANGUAGE: INSIGHTS FROM THE ODAM TILI THEORY. *World Scientific Research Journal*, 36(1), 143-145. Retrieved from <https://scientific-jl.com/wsrj/article/view/1904>
11. Odam Tili Akademiyasi. (n.d.). Core Principles. Retrieved from [Official OTA website URL].
12. Journal of New Century Innovations. (2025, August). THE ARCHAEOLOGY OF LANGUAGE. *Journal of New Century Innovations*, 83(1), 135-136. Retrieved from <https://scientific-jl.com/new/article/download/26674/25963/52053>
13. Magnus, M. (2001). What's in a Word? Studies in Phonosemantics [Doctoral dissertation]. Retrieved from <https://www.trismegistos.com/Dissertation/dissertation.pdf>
14. French, M. (2017). Cross-Linguistic Phonosemantics. University of Tennessee. Retrieved from https://trace.tennessee.edu/cgi/viewcontent.cgi?article=3099&context=utk_chanhonoproj
15. Fónagy, I. (2001). Languages within Language: An Evolutive Approach. John Benjamins Publishing. [As cited in French, 2017].
16. Aryani, A., Conrad, M., & Jacobs, A. M. (2013). Inherent emotional quality of human speech sounds. *Cognitive, Affective, & Behavioral Neuroscience*, 13(3), 506-517. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3620903/>
17. Capitol Technology University. (2023, May 30). The Ethical Considerations of Artificial Intelligence. Retrieved from <https://www.captechu.edu/blog/ethical-considerations-of-artificial-intelligence>
18. Abrahamson, E. (2023). The Coming Crisis of Complexity. *Harvard Business Review*. [Note: Placeholder citation for the concept; the actual source requires verification].