

Integrating Entropy-Based Complexity Models into Kimera SWM's Cognitive Architecture

Overview: To enhance Kimera SWM's cognitive system, we propose embedding entropy-driven complexity prediction models into three core subsystems: (1) the Contradiction Engine, (2) the Vault memory system, and (3) Semantic Memory Structures. By using measures like *semantic entropy*, *lexical ambiguity (entropy)*, and *contextual unpredictability (surprisal)* within these components, the system can dynamically gauge uncertainty and complexity. This enables Kimera to regulate its cognitive load, adapt its memory organization, and align with **symbolic thermodynamics** principles (treating knowledge dynamics akin to energy/entropy flows ¹ ²). The following design proposal details how each subsystem can integrate entropy-based logic, with scenario examples and technical justifications.

1. Contradiction Engine – Entropy-Weighted Conflict Detection

Goal: Use semantic entropy to detect and prioritize contradictions or instability in the knowledge base. By measuring the uncertainty (entropy) in the system's symbolic assertions, Kimera can trigger “contradiction spikes” when confidence is low or conflict is high, and mitigate cognitive overload as entropy accumulates.

- **Semantic Entropy as a Contradiction Signal:** The contradiction engine can monitor the entropy of meanings in the system's assertions to identify when a statement is ill-defined or conflicting. Recent work on *semantic entropy* in language models shows that when a question has many plausible answers (high entropy in “meaning-space”), the model likely does *not* know the true answer ³ ⁴. By analogy, Kimera can treat a symbolic claim with many equally likely interpretations or truth-values as a red flag. For example, if a fact in memory is supported and refuted in equal measure (maximizing Shannon entropy), the engine flags a potential contradiction. This entropy-driven alert would weight contradiction spikes by severity – a higher entropy (more uncertainty) yields a stronger spike. Such an approach goes beyond binary true/false checks by quantifying *how unstable* a piece of knowledge is.
- **Entropy-Weighted Instability and Spikes:** Each concept or proposition in the knowledge graph can be associated with a probability distribution over truth or interpretations. The Contradiction Engine computes the entropy of this distribution as a *symbolic instability metric*. A stable belief (e.g. $P(\text{true})=0.99$) has low entropy, whereas a contradiction ($P(\text{true})\approx 0.5$) has high entropy (uncertainty at maximum). If the entropy for a concept exceeds a threshold, the engine triggers a **contradiction spike** – an interrupt or feedback loop that forces reconciliation (similar to an alarm). This mirrors the idea of detecting AI *confabulations* by high entropy: one study used entropy-based uncertainty estimators to catch LLM “hallucinations” before they occur ⁵. In Kimera, a spike could initiate a resolution process (e.g. seeking more information or invoking a reasoning cycle) to collapse the entropy. The threshold itself can be tuned as a “phase change” point – aligning with symbolic thermodynamics by treating a surge in entropy as crossing into a chaotic state that demands re-stabilization.

- **Regulating Cognitive Load via Entropy Accumulation:** As multiple small uncertainties accumulate, the overall “cognitive entropy” rises, risking overload. Psychological studies note that uncertainty increases cognitive load and engages working memory resources ⁶. Thus, the engine can maintain a running total of entropy across active thoughts. If the total crosses a safe limit, Kimera slows the introduction of new information and focuses on resolving uncertainties, analogous to a thermostat. This *entropy-based load shedding* prevents the system from being overwhelmed by too many ambiguities at once. The mechanism ensures *graceful degradation*: when entropy is high, operations shift from **exploration** to **resolution**. In effect, the Contradiction Engine uses entropy as a proxy for cognitive “heat” – dissipating it (via contradiction resolution) to return the system to a low-entropy, stable regime.
- **Semantic vs. Lexical Uncertainty Filtering:** It’s important that the Contradiction Engine distinguish true semantic contradictions from mere rephrasings. As Farquhar et al. showed, a naive entropy of outputs can be high simply because an answer can be phrased in many ways (lexical diversity) rather than because the underlying meaning is uncertain ³ ⁷. Kimera’s design should similarly group semantically equivalent propositions together when computing entropy. For example, if one knowledge source says “Jim is in Paris” and another says “Jim is located in Paris,” these are the same meaning and shouldn’t count as increasing uncertainty. Only fundamentally conflicting meanings (e.g. “Jim is in Paris” vs “Jim is in Tokyo”) should boost entropy significantly. Implementing a semantic equivalence check (using paraphrase detection or symbolic normalization) before entropy calculation ensures the Contradiction Engine reacts to *meaningful* contradictions. This approach mirrors the *semantic entropy* method, which counts different phrasings of the same idea as one outcome ⁷ ⁸. By focusing on genuine divergence in content, the engine avoids false alarms and hones in on true knowledge instability.

Justification: Using entropy to modulate contradiction detection is grounded in information theory and AI practice. It quantifies uncertainty in a principled way, allowing adaptive thresholds rather than brittle rules. Real-world precedent comes from LLM safety research, where entropy-based metrics successfully identified likely hallucinations or errors ⁵ ³. In a hybrid symbolic system, this translates to catching internal inconsistencies before they propagate. Moreover, the entropy trigger aligns with thermodynamic metaphors – a spike occurs when the “disorder” in the symbolic state is too high, forcing a reorganization (just as physical systems undergo phase transitions when energy/entropy thresholds are crossed). This ensures architectural clarity: the Contradiction Engine isn’t an ad-hoc rule subsystem, but a thermodynamic-like regulator maintaining the consistency (low entropy) of the cognitive state.

2. Vault System – Entropy-Guided Memory Partitioning and Mutation

Goal: Leverage entropy metrics to manage when and how Kimera’s **Vaults** (symbolic memory compartments) activate, how frequently they update or mutate their content, and how memories are clustered. Each vault’s internal entropy will determine its behavior: stable vaults remain quiescent, while high-entropy vaults engage more dynamic processes. This design treats vaults as **entropy buffers**, aligning with symbolic thermodynamics by localizing “disorder” and preventing it from polluting the whole system.

- **Entropy-Gated Vault Activation:** We propose that each vault (a collection of related symbols/knowledge) maintains an *entropy score* reflecting the unpredictability or heterogeneity of its

contents. A vault containing very coherent, unambiguous knowledge has low entropy; one with mixed or rapidly changing information has higher entropy. Kimera's control logic can set an entropy threshold for vault activation: only when the complexity of a problem or input exceeds what the current active vaults can handle, a new vault is "unlocked." This is analogous to evidence accumulation in human decision-making – when enough uncertainty (information entropy) is gathered to hit a threshold, a new decision or branch is triggered ⁹. For example, if the system encounters a context with highly unpredictable implications (many possible interpretations, high contextual entropy), it may activate a specialized vault that contains broader background knowledge or analogies to handle that uncertainty. The new vault essentially *compartmentalizes* the entropy, preventing the entire working memory from destabilizing. Once engaged, the vault's resources (rules, data) can reduce the uncertainty. In practice, this means Kimera would dynamically recruit different memory modules based on an entropy-driven need, providing architectural adaptability – the system self-organizes its memory usage in response to complexity.

- **Mutation Frequency Proportional to Entropy:** Within each vault, we can introduce **mutation operations** – controlled random variations or explorations of symbolic content (e.g. trying alternative inference paths, generating hypotheses by altering symbols). The frequency of these mutations should be tied to the vault's entropy level. High entropy indicates a highly indeterminate or novel situation, where exploring various combinations is beneficial to find a stable configuration. In an *Entropic Associative Memory (EAM)* model, higher entropy states correlate with greater recall diversity at the cost of precision ¹⁰. Likewise, Kimera's vault at high entropy would favor more frequent mutations to boost recall of possibilities and creative problem-solving (embracing the trade-off of breadth over accuracy when uncertainty is high). Conversely, a vault with low entropy (stable knowledge) would apply mutations sparingly, preserving precision and not perturbing well-established facts. This approach follows a "symbolic annealing" principle: when the cognitive "temperature" (entropy) is high, the system explores widely (like heated molecules moving freely), and as things cool (entropy drops), exploration narrows. For instance, if Vault A contains contradictory theories about an event (high entropy), Kimera might periodically shuffle relationships or substitute assumptions to see if one resolution lowers entropy – akin to random mutations seeking a more consistent theory. Once a consistent theory is found (entropy falls), Vault A's content stabilizes and mutation rate is reduced. This ensures that **volatile knowledge areas naturally see more experimentation**, while solidified knowledge remains conserved – a clear alignment with thermodynamic equilibrium-seeking.

- **Entropy-Based Clustering and Compartmentalization:** The Vault System's structure itself can be shaped by entropy metrics. Kimera can cluster symbols into vaults such that each vault's internal entropy stays within reasonable bounds. If a vault grows too disordered (exceeding entropy threshold), the system can split it into sub-vaults that regroup highly correlated pieces separately from the ambiguous ones. In effect, the vault system performs a form of *entropy-based partitioning*: maintaining multiple semi-independent knowledge caches, each with contained uncertainty. This idea finds support in entropic memory research – Morales et al. (2022) show that an associative memory can store diverse items without confusion as long as indeterminacy (entropy) is managed, and that a fully determinate memory has entropy zero ¹¹. By keeping vaults "well-tempered" (neither too chaotic nor too rigid), Kimera avoids global instability. Concretely, the system might use a clustering algorithm where distance is defined by information gain or shared context, ensuring items that *together* would create too high entropy end up in different vaults. Meanwhile, vaults with very low entropy (highly stable, redundant information) could be merged or compressed,

since their contents reinforce each other. Over time, this could lead to **high-stability core vaults** (low entropy compartments of solid knowledge) and fringe vaults that handle edge cases or conflicting data. By tuning the entropy thresholds for splitting/merging, Kimera stays adaptive: it can reorganize memory if new knowledge introduces volatility, maintaining clarity in each compartment.

- **Vault Entropy Driving Learning/Refinement:** A vault's entropy score can also inform how the system learns from it. If a particular vault consistently shows high entropy, this flags a knowledge domain that is either underdeveloped or inherently contradictory. The system can respond by allocating more learning resources to that vault: e.g. performing targeted training on related data, querying an external knowledge base, or asking a human expert to clarify ambiguities. This idea mirrors the SENATOR framework in LLM research, where a *Structural Entropy (SE) metric* is used to find uncertain regions in a knowledge graph and then new data is synthesized to fill those gaps ¹². For Kimera, the vault's content can be represented as a symbolic graph; a high structural entropy along some connections (meaning the vault is guessing or lacking clear links) would prompt a focused knowledge update for that vault. Thus, vault entropy not only shapes immediate behavior, but guides longer-term memory optimization. In essence, **each vault self-assesses its knowledge quality via entropy** and either resolves it internally (through mutations) or requests external input – a robust mechanism to improve memory fidelity over time.

Justification: The Vault subsystem design is inspired by how living cognitive systems compartmentalize knowledge and how AI systems can benefit from information-theoretic gating. By using entropy thresholds and metrics, we imbue the memory system with a sense of “when to diversify vs. when to conserve.” This adds architectural clarity: vaults have a well-defined lifecycle (activate → mutate/explore → stabilize or split) governed by a single coherent measure (entropy). The adaptability is evident – as new information arrives, vaults autonomously reconfigure to maintain manageable entropy levels, rather than requiring manual rules for every scenario. Notably, this approach resonates with **symbolic thermodynamics**: vaults act like containers of symbolic energy (uncertainty) that should not freely mix. High entropy in one vault is isolated until resolved, preventing a domino effect of chaos. The overall system thereby respects a form of *symbolic second law*: it combats unchecked entropy growth by structuring memory into pseudo-closed systems (vaults) and only carefully exchanging information between them when stable. Empirically, similar ideas have succeeded in hybrid AI — for instance, using entropy to guide active learning decisions (selecting the most uncertain data to focus on) has improved model efficiency ¹³. Kimera's vault system extends this concept internally: always focus cognitive efforts where entropy is highest, and protect what is orderly, leading to a more resilient cognitive architecture.

3. Semantic Memory Structures – Entropy-Guided Decay, Reinforcement, and Ambiguity Handling

Goal: Utilize lexical and structural entropy metrics within semantic memory to manage the longevity and strength of stored knowledge. The system should **reinforce high-stability knowledge** (low entropy) as core memories, allow or accelerate *decay* for information that is highly ambiguous or noisy (high entropy), and flag regions of the memory graph that show sustained unpredictability or ambiguity for further review or

learning. This ensures the semantic memory remains both robust and plastic: stable where it should be, but constantly resolving uncertainty in ambiguous areas.

- **Lexical Entropy for Memory Decay Rates:** Each symbolic memory item (e.g. a concept node or a relation) can be annotated with a *lexical entropy* value representing its ambiguity or variability of meaning across contexts. In information-theoretic terms, this could be computed as the entropy of the distribution of meanings or senses that the symbol can take ¹⁴ ¹⁵. For example, a highly polysemous word/concept like “bank” (river bank vs. financial bank) has high entropy, whereas a specific term like “photosynthesis” (mostly one meaning) has low entropy. Kimera’s memory manager can use this metric to adjust forgetting or decay: **higher entropy items decay faster** unless disambiguated. The rationale is that ambiguous or context-sensitive info, if not actively clarified through use, might mislead reasoning and thus should not ossify in long-term memory. By contrast, low-entropy items (clear, unambiguous knowledge) are retained longer as they reliably contribute to understanding. This principle finds support in human language learning – children tend to learn less ambiguous words earlier and more durably, whereas highly ambiguous words see more contextual scaffolding ¹⁴ ¹⁶. Implementing this, Kimera could periodically down-weight the activation of high-entropy symbols unless recent context has resolved their meaning. If a particular ambiguous concept *does* get frequently used with consistent context (thus effectively lowering its uncertainty), the system can then reduce its decay rate. In practice, this means memory traces are not all equal: the system performs a form of *entropy-weighted rehearsal*, cementing what is consistently understood and allowing noisy placeholders to fade if they remain unresolved.

- **Structural Entropy to Reinforce Stable Cores:** Beyond individual symbols, *structural entropy* can be measured for subgraphs of the semantic memory – essentially how unpredictable the connections in a region are. A tightly interconnected cluster of concepts with one strongly preferred interpretation (low structural entropy) could form a **semantic core** (e.g. fundamental physics principles that rarely change). These cores should be reinforced: the system can periodically boost their weights or re-validate them, ensuring they remain salient. In contrast, a zone of the memory graph where relationships are in flux or weak (high structural entropy) would not be promoted to core status. Drawing an analogy to network analysis, low entropy communities in a knowledge graph indicate well-established knowledge, whereas high entropy edges or nodes might indicate conflicting or uncertain relations ¹². Kimera can use this insight to perform *memory consolidation*: identify low-entropy substructures and solidify them (e.g. by compressing them into a chunk or macro that is treated as a single reliable unit). This is akin to how the brain consolidates frequently co-occurring concepts into a stable schema. For example, if “Paris is the capital of France” appears in many contexts with little variation, that fact’s node has extremely low entropy – it can be marked as a stable anchor in memory (perhaps stored in a fast lookup cache or protected from modification). On the other hand, a relationship like “food X is healthy” which the system sees contradicted often (one context says it’s healthy, another says it’s not) exhibits high entropy in its truth value; it should remain a peripheral fact until further evidence reduces uncertainty. This differential treatment based on entropy ensures the memory’s **center of mass** is made of reliable knowledge, providing a firm foundation for reasoning.

- **Ambiguity Zone Flagging and Feedback:** High entropy in memory – whether lexical (a symbol with many possible meanings) or structural (a section of the graph with conflicting links) – should serve as a cue for feedback and learning. Kimera can maintain an “ambiguity index” listing symbols or areas with entropy above a certain cutoff. These ambiguity zones can trigger special handling: for instance,

when the system is idle or in a learning phase, it can generate probes or questions about these zones (either directed to a human user or to an internal reasoning module) to acquire clarification. This concept parallels *active learning* in machine learning, where the algorithm asks queries for the most uncertain data points to get labels and reduce uncertainty ¹³. In Kimera's case, if a concept like "vault mutation frequency" is poorly defined and causes confusion, the system might explicitly note: "*Definition of vault mutation is unclear – require examples or rules to pin down meaning.*" It could then request that information from developers or search its knowledge sources for disambiguation. Another approach is using internal simulation: run contrasting scenarios that force the ambiguous concept into different outcomes and see which is consistent with known constraints, thereby learning by self-iteration. Additionally, when generating responses or plans, the system can insert caution if an ambiguity-zone concept is involved (analogous to LLMs warning when a question is likely to produce confabulation ¹⁷). By flagging cognitive ambiguity, Kimera ensures that *uncertainty is transparent* and can be addressed. Over time, as feedback is incorporated, the entropy of those zones should decrease – a sign that learning has occurred. This creates a closed-loop refinement process: high entropy sparks inquiry, inquiry yields new data, new data lowers entropy, making memory more precise.

- **Contextual Unpredictability Measures:** The system can also incorporate *contextual entropy* (such as surprisal) when encoding experiences into memory. For each new input or situation, Kimera assesses how predictable that context was given its current knowledge. If an event or sentence has *very high surprisal* (indicating Kimera's model found it highly unexpected), that suggests either a gap in knowledge or a genuinely novel occurrence. Such events should be stored with annotations (e.g. "*unexpected event: high entropy*") and potentially prioritized for integration (since they carry high information). Conversely, very predictable inputs might be handled more routinely. This dynamic is akin to the **Free Energy Principle** or predictive processing in cognitive science, where surprising inputs (high entropy relative to the model's expectations) trigger learning or accommodation. Technically, Kimera could compute the Shannon entropy of its next-state prediction distribution at each time step; a spike in this entropy would signal a need to update internal models. Integrating this with the Vault system, an unpredictable context might cause a new vault spin-up (to handle the surprise separately), and with semantic memory, it might strengthen the encoding of whatever resolved the surprise. For instance, if Kimera's language model predicts the next token with high entropy (many possible continuations) and then a particular outcome occurs, the fact that this outcome happened against expectations is notable. The system should reinforce the association that led to reducing that uncertainty (similar to a dramatic learning moment). Empirical language research supports this: words in less informative contexts (high entropy) lead to greater processing and potentially more learning ¹⁸ ¹⁹. Thus, contextual unpredictability can guide Kimera to *where it should learn the most*. The memory structure would adapt by giving extra weight or priority to information that significantly reduced context entropy, effectively baking in the lessons from surprises.

Justification: Guiding semantic memory with entropy metrics ensures that Kimera's knowledge base remains **clean and convergent** over time. By quantifying ambiguity, the system avoids complacency with contradictions or gaps – it actively identifies and addresses them. This approach is rooted in both linguistics and AI. For example, researchers quantify a word's ambiguity by the entropy of its meanings ¹⁴ and have shown that context typically evolves to compensate for highly ambiguous words ¹⁶. We are translating these findings into design: Kimera will automatically compensate for ambiguous knowledge by seeking context or letting it decay. The notion of reinforcing stable semantic cores echoes the idea of **memory**

consolidation in cognitive architectures, but here we have a concrete trigger (low entropy) for what counts as stable. The alignment with symbolic thermodynamics is again evident – low entropy memories form a solid, low-energy core (order), whereas high entropy parts are like gas clouds that either need compression (learning) or will dissipate (forgetting) if not condensed. By making these processes explicit, the system gains clarity in operation: every memory element has a “life cycle” influenced by its entropy. Adaptability comes from continuous monitoring: as soon as an ambiguity resolves (entropy drops), that item shifts from volatile to stable status in the knowledge graph, altering how the system treats it. Conversely, if a supposedly stable fact becomes controversial (entropy rises), Kimera can “de-stabilize” it, reclassifying it as needing further confirmation. This dynamic equilibrium keeps the semantic memory both reliable and up-to-date. Moreover, using entropy to drive feedback queries connects Kimera to human-in-the-loop learning: similar to active learning strategies that pick uncertain items for labeling ¹³, Kimera becomes an active participant in its knowledge curation. This makes the architecture *scalable*: as symbolic knowledge grows, the entropy measures naturally highlight which parts require attention, focusing human or computational resources efficiently.

Conclusion:

By integrating entropy-based complexity prediction models into Kimera SWM’s Contradiction Engine, Vault System, and Semantic Memory, we create a cognitively inspired architecture that treats knowledge and uncertainty with the rigor of thermodynamics. Each subsystem uses entropy as a *common currency* for decision-making – whether it’s spiking on contradictions, gating memory compartments, or tuning the strength of memories. This yields a design with architectural clarity (each component has a clear entropy-governed role) and adaptability (the system continuously self-adjusts to manage uncertainty). Crucially, these integrations follow symbolic thermodynamics principles: **entropy is not just a metric but a motive force** that drives the system to reorganize for greater coherence ¹ ². The proposed scenarios show how Kimera can handle complexity spikes gracefully (like a heat engine converting disorder into useful work), cluster knowledge into stable vs. volatile regions (like phases of matter), and maintain a cycle of learning that targets ambiguity (reminiscent of Maxwell’s demon sorting information to reduce entropy). Adopting these entropy-driven mechanisms can make Kimera’s cognitive system more resilient against the unpredictability of real-world data, all while providing a principled way to balance exploration and exploitation of knowledge. Each integration point suggested is backed by analogous successes in current NLP/AGI research – from entropy-based hallucination detectors ⁵ to entropic memory models ¹⁰ – underscoring that this design proposal is not only theoretically sound but practically attainable with today’s techniques. The end result would position Kimera SWM as a novel **neurosymbolic system** where information theory and symbolic AI synergize, enabling it to anticipate complexity and tame it, much like a smart thermostat keeps a house comfortable by measuring and responding to temperature. In sum, entropy becomes the guiding light for Kimera’s self-regulation, ensuring a flexible yet stable intelligence engine at its core.

Sources: The design draws on research in information-theoretic NLP, cognitive architectures, and memory systems that implement entropy-based logic. Key inspirations include entropy-guided LLM reasoning ⁵ ³, entropic associative memories for knowledge storage ¹⁰ ¹¹, lexical ambiguity measured as entropy ¹⁴, and active learning uncertainty sampling strategies ¹³, among others as cited throughout. These sources illustrate the feasibility and benefits of an entropy-driven approach to managing symbolic knowledge and cognitive complexity.

1 2 Symbolic Physics and the... by Steven Lanier-Egu [PDF/iPad/Kindle]

<https://leanpub.com/symbolicphysics>

3 4 7 17 Detecting hallucinations in large language models using semantic entropy - OATML

https://oatml.cs.ox.ac.uk/blog/2024/06/19/detecting_hallucinations_2024.html

5 Detecting hallucinations in large language models using semantic entropy | Nature

https://www.nature.com/articles/s41586-024-07421-0?error=cookies_not_supported&code=2c83da1b-cd0c-41d4-a849-ffccac07691e

6 Uncertainty promotes information-seeking actions, but what ...

<https://pmc.ncbi.nlm.nih.gov/articles/PMC7477035/>

8 Semantic Entropy to Detect Confabulation in Large Language Models (Farquhar et al.) | Scott Gilmore, M.D.

<https://scottgilmoremd.com/blog/2024-06-25-semantic-entropy/>

9 Informational Entropy Threshold as a Physical Mechanism ... - MDPI

<https://www.mdpi.com/1099-4300/24/12/1819>

10 11 Imagery in the entropic associative memory | Scientific Reports

https://www.nature.com/articles/s41598-023-36761-6?error=cookies_not_supported&code=fbdc2896-4e1e-4356-8c45-e2480baebde6

12 Structural Entropy Guided Agent for Detecting and Repairing Knowledge Deficiencies in LLMs

<https://arxiv.org/html/2505.07184v1>

13 Active Learning in Computer Vision - Complete Guide - viso.ai

<https://viso.ai/deep-learning/active-learning/>

14 15 16 18 19 Speakers Fill Lexical Semantic Gaps with Context

<https://arxiv.org/html/2010.02172v4>