

Kimera SWM System Architecture

Introduction

Kimera SWM (Semantic Working Memory) is a novel cognitive architecture for managing knowledge and context in an AI system. It employs an **entropy-guided semantic memory system** with mechanisms for contradiction detection, stable knowledge storage ("vault"), and a **non-token-based knowledge representation**. Memories are stored as event-based units called "**scars**", reflecting lasting traces of experiences or inputs. The architecture is strictly layered – comprising **Input/Ingestion**, **Processing/Core Logic**, **Semantic Storage**, and **Interaction/Exchange** layers – to ensure clear modular separation. Internal data flows and memory operations are guided by principles analogous to thermodynamic entropy, aiming for efficient information preservation and minimal redundant energy usage. Importantly, the Kimera SWM core is a zero-dependency kernel: it relies on no external frameworks within its core logic, ensuring full control over data structures and operations. This document outlines the system architecture in a formal engineering format, detailing layer responsibilities, data flows, key data structures, interface definitions, and design constraints. Emphasis is placed on clarity, adaptability for future evolution, and testability of the system.

Architecture Overview

The Kimera SWM architecture is organized into four primary layers, each with distinct roles and well-defined interfaces:

- **Input/Ingestion Layer:** Handles incoming data or events from external sources. It translates raw inputs into the internal semantic format (using the EchoForm grammar) for further processing.
- **Processing/Core Logic Layer:** The "brain" of the system containing core modules such as the **Contradiction Engine**, **Entropy Manager**, and **Causal Propagation logic**. This layer performs analysis, consistency checks, and routes information between input and storage, following entropy-guided rules.
- **Semantic Storage Layer:** Provides persistent semantic memory. It includes a **Memory Vault** for long-term, structured knowledge and a **Scars Repository** for event-based memory units. This layer stores information in a non-token-based representation for efficient retrieval and update.
- **Interaction/Exchange Layer:** Defines how external entities (users, applications, other systems) interact with Kimera SWM. It offers interfaces (APIs) for querying the memory system and exchanging information, formatting inputs/outputs as needed without exposing internal complexities.

Figure: Kimera SWM high-level architecture overview. The system is organized into layered modules: External sources provide data to the Input/Ingestion layer, which parses inputs via the EchoForm grammar. The Core Logic layer processes inputs, checks for contradictions, and manages entropy-based decisions, interacting with the Semantic Storage layer (Memory Vault and Scars repository). External clients interact via the Exchange layer's query interface. Dashed arrows indicate internal checks/notifications (e.g., contradiction checks, causal propagation), and dotted arrows indicate retrieval operations.

Each layer communicates with adjacent layers through clearly defined interfaces, ensuring modular independence. For example, the Input layer passes normalized events to the Core logic via an internal event object, and the Core uses storage via a memory interface abstraction – neither layer depends on the internal implementations of the others. This design facilitates **extensibility** (each module can be modified or replaced independently) and **testability** (modules can be unit-tested in isolation using mock interfaces).

Design Principles and Constraints

Several key design principles underlie the architecture:

- **Entropy-Guided Memory Management:** Decisions about storing, recalling, or pruning information are based on information-theoretic measures of “semantic entropy” rather than arbitrary heuristics. High-entropy (high-information or novel) content is prioritized for retention, while low-entropy or redundant content may be safely pruned to optimize capacity. This aligns with thermodynamic principles like Landauer’s limit, treating information deletion as a costly operation to minimize ¹.
- **Consistency via Contradiction Handling:** The system maintains a coherent knowledge base by detecting and handling contradictions. New information that conflicts with stored beliefs triggers a **belief revision** process similar to a truth-maintenance system – identifying and resolving conflicting beliefs so the knowledge base remains internally consistent ².
- **Non-Token-Based Representation:** Internally, knowledge is represented in a conceptual or vectorized form rather than as raw text tokens. This allows semantic content to be stored abstractly (for example, as nodes and relationships in a graph, or as high-dimensional embeddings) without direct reliance on language tokens. It enables language-agnostic reasoning and compact storage of meaning.
- **Event-Based Memory Units:** Each significant input or occurrence is encapsulated as an **Event “Scar”** – an immutable record with semantic content and context metadata (timestamp, source, links to related events). Scars serve as atomic memory entries that can be composed into larger narratives or knowledge structures. The term “scar” highlights that each event leaves a lasting mark on the system’s memory state.
- **No External Dependencies in Core:** The core logic layer is implemented with a zero-dependency philosophy – avoiding external libraries or frameworks. All essential data handling (parsing, storage management, entropy calculation, etc.) is done with custom or standard library code. This ensures determinism, security, and full transparency of the system’s inner workings, which is critical for debugging and for extending the system without being tied to third-party changes.
- **Extensible Interfaces:** Each layer exposes interface points (APIs, data exchange formats) that can be extended. For instance, new input handlers can be added to the ingestion layer for different data types, or additional query types can be supported in the interaction layer, without altering the core logic. This modularity supports ongoing evolution of the system.

With these principles in mind, we now detail each architecture layer, followed by the dynamic behaviors (entropy management, memory activation, causal propagation) that govern the internal data flow.

Layered Architecture

Input/Ingestion Layer

Role: The Input/Ingestion layer is responsible for capturing and normalizing all incoming information to the system. Inputs may include user-provided text, sensor readings, events from other programs, etc. This layer ensures that raw data is translated into the internal semantic format for the core logic to process.

Components:

- **Input Adapters:** Modules or APIs that interface with various external sources (e.g., text streams, file inputs, sensor feeds). Each adapter handles communication protocols or data formats of the source, abstracting them into a common format. - **EchoForm Parser:** The parser uses the **EchoForm grammar** (the system's orthographic/semantic substrate) to convert raw input into a structured internal representation. EchoForm serves as a formal language in which inputs are represented as semantic structures (similar to parse trees or logical forms). For example, a sentence or an event description is parsed into an EchoForm syntax tree capturing entities, relationships, and context. This parser ensures that no matter the input source, the core logic always receives semantically organized data.

Data Flow: When an input arrives, the Input Adapter first captures it (for instance, reading a JSON event or a user query string). The raw content is then fed into the EchoForm Parser, which produces a **Parsed Event Object** – essentially an `EventScar` data structure populated with the semantic content of that input. This object includes fields such as a unique event ID, timestamp, the semantic content (in EchoForm form or an equivalent object graph), and initial metadata (source, priority, etc.). The parsed event is then forwarded to the Processing/Core layer via a defined interface call (e.g., `Core.ingestEvent(event)`).

Interface & Boundary: The Input layer's interface to the Core is a method or message queue that accepts event objects. The Input layer does not need to know how the Core uses the data; it simply hands off a fully parsed event. Conversely, the Core logic does not need to know details of how parsing was done or what the raw format was – it trusts that the incoming object adheres to the EchoForm-defined structure. This separation means new input formats can be supported by adding a new parser module without any change to the core processing.

Example: If an external user inputs: *"The skyship detected a storm at 3PM."*, an Input Adapter for user text receives this string and passes it to the EchoForm Parser. The parser (using EchoForm grammar rules) might produce an internal representation such as `(DETECT EVENT)[subject: skyship][object: storm][time: 1500hrs]`. This structured event (as a new EventScar instance) is then sent into the Core logic.

Processing/Core Logic Layer

Role: The Processing/Core Logic layer is the central processing unit of Kimera SWM. It accepts structured events from the Input layer, performs analysis and integration of that information, and updates or queries the Semantic Storage accordingly. The Core layer is also responsible for reasoning tasks like contradiction checking, entropy calculation, and managing how memories are connected or activated. It outputs results or retrieved information to the Interaction layer upon request. All major cognitive functions are implemented here with no external dependencies.

Subcomponents: The Core layer is subdivided into several key modules that work together:

- **Semantic Processor:** A central coordinator that handles incoming events and queries. It interprets EchoForm-structured data and decides which subsystems to engage. For a new input event, the Semantic Processor orchestrates the sequence: (1) analyze the event content, (2) consult the Contradiction Engine for consistency, (3) use the Entropy Manager to decide storage, and (4) update the memory structures (Vault/Scars). For a query, it triggers a memory search and composes a response (possibly using EchoForm grammar in reverse to form a natural language answer). The Semantic Processor essentially implements the core logic flowchart, ensuring the correct order of operations and combining results from various subsystems.
- **Contradiction Engine:** This module ensures **knowledge consistency**. When new knowledge enters, the Contradiction Engine checks it against existing beliefs in the Memory Vault. It employs logical inference rules (or potentially semantic embedding similarity) to detect direct contradictions or violations of constraints. If a contradiction is found, the engine triggers a resolution policy. By default, the policy may mark the new information as suspect or create a **contradiction record** linking the conflicting facts. More advanced resolution might involve **belief revision** – e.g., removing or downgrading the confidence of one of the conflicting beliefs to restore consistency ². The Contradiction Engine ensures the system is aware of any internal conflicts. It could be implemented akin to a Truth Maintenance System, keeping justification graphs of facts and updating beliefs whenever inputs are retracted or added to maintain coherence.
- **Entropy Manager:** The Entropy Manager governs decisions related to information preservation and pruning, guided by **semantic thermodynamic entropy** principles. Each incoming event or stored memory is evaluated for its information content and relevance. The Entropy Manager may compute an entropy score for an event (for example, based on its surprise or how much it decreases uncertainty in the knowledge base). This score influences whether the event is stored prominently in the Memory Vault, stored only as a low-priority scar, or discarded. Over time, the Entropy Manager also monitors the **entropy of the memory system** – for instance, if certain memory areas become too chaotic (high entropy) or stale (low entropy), it may initiate pruning or compression. **Entropy-governed pruning** means that the system preferentially removes or archives information that provides minimal semantic utility (e.g., repeated facts or low-impact details) while **preserving meaningful content** ³ ⁴. The Entropy Manager works like a garbage collector tuned by information theory: it might use thresholds such that any memory contributing below a certain mutual information value to the rest of the system is pruned to conserve resources. Importantly, outright deletion is avoided unless necessary – information may be compressed or moved to long-term archive instead, echoing Landauer’s principle that erasing information has a minimum energy (or cost) ¹. In practice, this could mean marking seldom-used scars as “inactive” (non-recallable) rather than deleting them, or merging similar memory records.
- **Causal Propagation (Mirror Logic):** This subsystem manages the **cause-effect relationships** between events and knowledge, using a *mirror-based propagation* approach. Every time a new event (scar) is stored, the Causal Propagation module updates links between the cause and effect represented. For example, if Event B is recorded as an outcome of Event A, the system creates a bidirectional association: $A \rightarrow B$ (cause leads to effect) and a *mirror* link $B \rightarrow A$ indicating effect can be traced back to cause. The propagation logic ensures that if a cause has multiple effects or an effect has multiple potential causes, these relationships are captured in a graph structure. It then *propagates* any necessary updates: e.g., if a cause event’s status changes (say it is invalidated by the contradiction engine or pruned), the mirror logic can propagate this change to flag downstream effects as potentially invalid or in need of review. Similarly, if a new event reinforces a causal chain,

the mirror logic might strengthen those links. This mechanism essentially maintains a **causal graph** of the knowledge base, enabling reasoning like “if A happened, then later B happened” and also the reverse “B happened, possibly because of A”. The “mirror” aspect implies symmetry in how cause/effect information is stored for robustness of retrieval (any direction can be traversed). This provides a structural way for the system to understand temporal and causal dependencies among scars.

- **Non-Token Semantic Representation:** Although not a separate module, it's worth noting that within all core processing, data is handled in a semantic form (often graph or vector form) rather than raw text. The Semantic Processor and related algorithms manipulate concepts, relations, and embeddings. For example, a concept like “storm” is represented by a semantic identifier or embedding that may connect to attributes (like location, time) and related concepts (like “weather” or “skyship detection”). This allows the core to perform reasoning (contradiction detection, similarity, activation spreading) on a conceptual level. The EchoForm grammar plays a role here by providing a structured **orthography** for these concepts, but once parsed, the core deals with an object model or graph of the concept, not with strings of words.

Internal Interactions: Within the Core layer, these subcomponents interact through shared data structures and events: - The **Semantic Processor** invokes the Contradiction Engine by providing it the new event's assertions and asking for a consistency check against the Memory Vault. The output might be a simple boolean (conflict found or not) plus details of what conflicts. - The **Entropy Manager** is invoked to evaluate the new event (and possibly the overall state). It may attach an entropy score to the event and decide on its fate (e.g., “accept into vault”, “store only as scar”, or “discard”). It could also trigger a prune cycle at certain intervals. - If the event is accepted, the **Semantic Processor** then writes it into the Semantic Storage: it creates a new entry in the Scars Repository (for historical record) and potentially updates or adds to the Memory Vault (for structured recallable knowledge). If the event is judged low priority, maybe it stays only in Scars (which might be a log), whereas significant structured info is merged into the Vault's knowledge base. - Upon storing a new event, the **Causal Propagation** module is notified (perhaps via an internal event or by the Semantic Processor). The module checks if this new event has references to prior events or if it can be linked as a cause/effect. It then updates the internal causal index accordingly. This may involve reading related scars from the repository or vault (for example, finding if any existing event is marked as a predecessor of this type of event) and then linking them. These operations do not require external input – it's purely internal book-keeping to enrich the knowledge graph.

Interface & Boundaries: The Core Logic provides interfaces both “northbound” (to the Interaction layer for queries) and “southbound” (to the Storage layer for memory access). However, these interfaces are abstract: - To Storage: Instead of directly manipulating database records, the Core uses a **Memory Storage Interface** (e.g., an object with methods `saveScar(event)`, `updateVault(concept)`, `find(query)` etc.). This abstracts whether the storage is in-memory, file-based, or a database. It respects the zero-dependency rule by using either custom data structures or minimal standard features (e.g., file I/O) rather than an external DB engine. - To Interaction: The core exposes a **Query Interface** that the Interaction layer calls. This could be a function like `answerQuery(query_object)` or a message queue that the Interaction layer submits queries to. The Core then processes the query (often by searching the Memory Vault and Scars) and returns a result.

Importantly, the Core does *not* depend on how the query came in (could be via an API call, CLI, etc.) nor on how the results will be used – it simply receives structured queries and returns structured answers (which the Interaction layer can format for the user). Similarly, the Core doesn't know the details of how data is physically stored; it relies on the storage interface to handle actual read/write operations.

We will further explore the dynamic behaviors (like memory activation and pruning) in a later section, but first we describe the storage structures themselves.

Semantic Storage Layer

Role: The Semantic Storage layer is responsible for maintaining the system's knowledge over time. It holds both the long-term consolidated knowledge (in the **Memory Vault**) and the chronological log of experiences (the **Event Scars Repository**). This separation allows the system to differentiate between core, structured knowledge and raw historical data. The storage layer ensures efficient retrieval, update, and pruning of memory items in service of the Core's needs.

Components:

- **Memory Vault:** A structured knowledge base that contains distilled semantic information. The vault can be thought of as a curated store of facts, concepts, and stable relationships that the system has accepted. It is organized in a way that supports semantic queries – for example, as a network/graph of nodes (concepts or entities) and edges (relationships), or as a set of logical assertions/triples. The vault likely includes indexes or maps by concept name or identifier for quick lookup. It may also maintain truth values or confidence levels for each entry (especially if contradictions have occurred, some facts might be marked as uncertain or superseded). The vault is **hierarchical**: knowledge could be organized in categories or ontologies such that general concepts are linked to more specific ones, forming a tree or graph. This hierarchy is leveraged for **“glowing” memory patterns**, described later, where activation can spread through these links. - **Event “Scars” Repository:** A chronological store of all significant events or inputs the system has encountered. Each entry here is an **EventScar** object, which includes the details of that event (in EchoForm structured content) as well as metadata such as timestamp, source, related links, and entropy score at creation. The scars repository can be imagined as an append-only log of events (much like an event sourcing system in software architecture). It preserves the history, which is useful for auditing, temporal reasoning, or retraining the vault if needed. Scars may also contain pointers to related scars (forming a graph of events) or to corresponding entries in the Memory Vault (if the event led to knowledge in the vault). Because scars are time-ordered, they also facilitate the causal propagation logic – e.g., by looking at sequences of events. - **Indexing and Lookup Structures:** In addition to the main data stores, the storage layer can maintain indices for efficient search. This could include a semantic index (mapping keywords or concept embeddings to vault entries), a temporal index for scars (to quickly retrieve events in a time range or sequence), and a contradiction index (linking contradictory entries for quick reference by the contradiction engine). These are implementation details, but critical for performance: we need to be able to find relevant knowledge quickly when a query comes or when checking a new input.

Data Structures: The storage defines clear formats for its contents. For example, an **EventScar structure** might be defined as:

```
EventScar {
  id: UUID or incremental ID
  timestamp: DateTime
  content: EchoForm structured content (semantic graph or AST of event)
  links: [IDs of related EventScars] (e.g., predecessor, cause, similar event
  references)
  derived_vault_entries: [IDs of Vault entries] (if this event resulted in
```

```
new/updated knowledge in vault)
  entropy_score: float (information content measure at insertion)
  status: {active, archived, conflict_flag} (flags for pruning or
contradiction resolution)
}
```

And a **Vault entry structure** could be something like:

```
VaultEntry {
  key: Concept or Assertion ID (e.g., "Skyship -> detects -> Storm")
  content: semantic content (graph node or logical assertion)
  confidence: float (belief strength or validity)
  links: [related VaultEntry IDs] (e.g., hierarchical parent/child or
associative links)
  source_events: [IDs of EventScars] (events that introduced or support this
knowledge)
  last_accessed: timestamp (for usage-based pruning)
}
```

(These are conceptual schemas to illustrate how data might be organized.)

Interactions: The Storage layer primarily interacts with the Core logic layer: - When the Core needs to store a new event, it calls into the repository (e.g., `ScarsRepo.add(event)`), which appends the scar and returns the assigned ID. If the event is to be integrated into long-term memory, the Core then calls `Vault.upsert(concept)` or similar, to add a new vault entry or update an existing one (for example, increasing confidence if it was repeated). - When the Core (Semantic Processor) needs to answer a query, it will call `Vault.find(query)` to retrieve matching knowledge. Under the hood, this might use a semantic index or perform a graph traversal. If nothing is found in the vault, the system might also search recent scars via `ScarsRepo.search(criteria)` especially for very recent context not yet integrated. - The Contradiction Engine can query the vault for any fact or concept that matches an incoming event's content to check for negation or inconsistency (e.g., if new event says "X is true", check if vault has "X is false"). This is done through structured queries to the vault (logical queries). - The Entropy Manager might iterate through vault entries or scars periodically (or on demand) to identify candidates for pruning. For example, it might call something like `Vault.listEntries()` with a filter for low usage or low entropy, then remove them or mark them archived. Similarly, old scars might be archived: the repository could move them to an archive file or mark them inactive if their timestamp is beyond a retention window and they have low significance.

Constraints: The storage is designed to be efficient and consistent: - It should support ACID-like consistency for updates (especially for the vault) to ensure that if the system is interrupted, the knowledge base isn't left in a contradictory half-updated state. This may involve simple transactional logic if built in-house. - Since no external database is used in the core, the storage may rely on in-memory structures with periodic snapshotting to disk, or use lightweight file storage (like append-only log for scars, and a serialized graph for vault) when persistence is needed. This is a trade-off to avoid external dependencies: the system must manage persistence itself. - The data structures need to be optimized for the typical access patterns: frequent reads for query (so indexing is key), and writes on each new event (which must be fast, likely using

append operations). - Scalability is a consideration: as the number of scars grows, the system might implement a **vault compaction** process where old scars that are fully assimilated into the vault and no longer independently useful can be removed or offloaded to a long-term archive outside the active memory.

Interaction/Exchange Layer

Role: The Interaction/Exchange layer handles all communication between Kimera SWM and the outside world. It provides the user interface or API endpoints through which external users or systems can send queries, receive answers, or subscribe to updates. In essence, this is the **presentation layer** of the architecture, translating between the internal representations and user-facing formats.

Components:

- **Query Interface (API):** A defined interface (which could be a function call in a local environment, a RESTful API, a CLI, etc., depending on deployment) that accepts queries or requests. A query could be a question like "What storms were detected by the skyship?" or a request like "Explain the last event" or even a command "Integrate this new data...". The Query Interface ensures the incoming request is parsed (possibly using EchoForm if the query is in natural language) or interpreted into the internal query format. It then calls the Core's query answering function and formats the result for output. - **Response Formatter:** After the core produces an answer (likely as structured data or as an EchoForm representation of the answer), the Interaction layer formats it into the desired output form. If the client expects a natural language answer, this component might use the EchoForm grammar in reverse to generate a grammatical sentence or a JSON if a structured response is needed. It might also attach any metadata (like confidence or references) if required by the application. - **Exchange Protocols:** If Kimera SWM is part of a larger system, the exchange layer might also handle subscription to certain events or pushing updates. For example, if an external system wants to be notified whenever a contradiction is detected in the knowledge base, the Interaction layer could provide a callback mechanism or event feed for that. Or it might handle multi-turn interactions, maintaining a session context with a user. - **Security/Governance Filters:** (Optional) This layer can also enforce any access control or filtering. For instance, certain vault information might be classified, and the Interaction layer can ensure those are not exposed via the API unless proper authorization is present. It can also throttle requests or sanitize inputs to protect the core.

Data Flow: When an external client makes a query, the flow is: 1. The query arrives at the Query Interface (say via an HTTP request or function call). The query may be in natural language or a structured format. If natural language, the EchoForm Parser (or a similar query parser) is used here to interpret it into an internal query object. For example, "Has the skyship detected any storms?" becomes a structured query asking for events of type detection by agent "skyship" where object is "storm". 2. The Interaction layer calls the Core logic through the defined Query API, passing this query object. This is a clean call like `Core.answerQuery(queryObject)`. 3. The Core (Semantic Processor) executes the query: it searches the Memory Vault (and possibly Scars) for matching information. Suppose the vault has an entry like `Skyship -> detected -> Storm at LocationX on TimeY`, the core finds this and compiles an answer. The answer might be formulated internally as another EchoForm structure or a simple data record: e.g., an answer object containing the fact(s) found or a statement like `(YES)[Skyship detected storm at 3PM]`. 4. The Core returns the result to the Interaction layer. The Response Formatter then takes this and converts it into the client's expected format. If the client asked in English, it might produce a sentence: "Yes – the skyship detected a storm on record (on 3PM at location X)." If the query was via API expecting JSON, it might output a JSON with fields for the event details. 5. The response is then sent back to the external client.

In a request-response scenario, that's the end. In a continuous interaction scenario, the Interaction layer might then wait for the next query or prompt from the user, possibly keeping some state if needed (though ideally the core is stateless aside from its memory, any conversational context could also be stored as part of the memory).

For outgoing communications like event notifications, the Interaction layer would work in reverse: an internal trigger (say the contradiction engine flags a conflict) would be picked up by a notifier component which then formats an external message or signal.

Modular Boundaries: The Interaction layer is designed such that it *only* deals with input/output formatting and protocol handling. It does not perform any deep logic on the content – that is the core's job. This means the system could have multiple frontends (e.g., a web UI, a command-line tool, a direct function call interface) all utilizing the same core logic underneath. Each frontend would just translate its input into the core's query format and convert core's answers to its output medium. This separation allows, for example, the core to be tested directly with known query objects, bypassing the external interface.

Adaptability: Thanks to this design, adding a new mode of interaction (for instance, a voice interface or integration with a chat platform) would only require changes in the Interaction layer – mapping spoken input to text/EchoForm and reading out answers – without touching core logic or storage. This fulfills the extensibility goal and makes the system more versatile.

Internal Dataflows and Semantic Entropy Management

This section describes how data moves through the system in typical scenarios and how **semantic thermodynamic entropy principles** inform those flows. We consider two primary scenarios: **Knowledge Ingestion** (a new event coming in) and **Query/Recall** (answering a question or retrieving info). We also detail how **non-recall, memory activation, and entropy-based pruning** are applied during these processes.

Ingestion Flow (Event Processing)

When a new event or piece of information is ingested, the following steps occur (with references to components in previous sections):

1. **Capture and Parse:** The Input layer's adapter captures the raw event and invokes the EchoForm Parser, yielding a structured EventScar object for the new information. (For example, raw text -> parsed semantic representation).
2. **Core Ingestion & Preliminary Analysis:** The Semantic Processor in the Core receives the EventScar. It examines the content and metadata. Immediately, it may tag the event with some context or priority (for instance, if the event source is high-trust or the type of event suggests it's critical).
3. **Contradiction Check:** The Semantic Processor passes the event to the Contradiction Engine to verify consistency. The engine queries the Memory Vault for any entry that logically conflicts. If, say, the event asserts a fact that an existing vault entry negates, a conflict is found. In that case, the Contradiction Engine might annotate the event (or the existing entry) with a contradiction flag. System policy may be to prefer existing knowledge over new (requiring new info to be corroborated before acceptance), or to accept new info and mark old as dubious – this policy is configurable but

managed within this engine. Regardless, the contradiction check ensures the system *knows* about the inconsistency. If no direct conflict is found, the event is considered consistent and can proceed.

4. **Entropy Evaluation:** The event is then evaluated by the Entropy Manager. The manager computes an **entropy value** for the event – essentially measuring how much new information it carries relative to what’s already known. This could involve checking how many new concepts or relations it introduces vs. how expected it was. For example, an event identical to many past events has low information (high redundancy), whereas a novel event type has high entropy (surprising information). Based on this value and possibly the current “entropy budget” of the system (how much free capacity or complexity it can handle), the manager makes a decision:
5. If the event’s entropy is above a threshold (very informative) or it’s otherwise deemed important (perhaps by source or type), it is marked for full retention (store in vault and scars).
6. If moderate, maybe only store in the scars (it’s recorded but not elevated to structured knowledge in vault).
7. If very low (e.g., a duplicate or something already known), the system might choose not to store it at all (or just increment a counter that this info was repeated). This is the **non-recall** in action at storage time – not every input is remembered.
8. The entropy evaluation might also trigger **pruning** in parallel: if we added something new and memory is reaching capacity, the manager might find an old low-entropy scar to remove or compress, keeping total entropy in balance (like a cache eviction based on info content).
9. **Storage Update:** Given the decisions from above, the Semantic Processor proceeds to update storage:
10. The EventScar is added to the Scars Repository regardless (unless it was wholly rejected). If it was deemed not worth keeping, it might still log minimal info like “event received at time X but skipped due to low entropy” for audit, then not store details.
11. If the event carries new factual knowledge or updates, the Memory Vault is updated. This could mean adding a new VaultEntry (for a brand-new fact or concept) or updating an existing entry (e.g., increasing its confidence, adding a new relation). The vault update is done through its interface ensuring consistency (for example, linking the new entry to existing related entries, updating indices).
12. The event’s scar record may store the ID of any new/updated vault entries (for traceability).
13. **Causal Link Propagation:** After storing, the system registers any causal relations. If the event explicitly has a causal context (say the input indicated “because A happened, B happened”), those are stored. If not explicit, the system may infer some sequence: for instance, if an event has a timestamp and there was a known prior event from the same source or context, it might link them as sequential (though causality is not guaranteed by sequence alone, the system could note a possible chain). The Causal Propagation module updates the cause-effect graph with this event. It will mark, for example, “Event #123 follows Event #118” and if a cause relation is known (“triggers” or “results in”), it will mirror-link them. This step ensures the timeline and dependency information is captured.
14. **Post-ingestion Hooks:** With the new event stored, any registered listeners or triggers are notified. For example, if the event created a contradiction flag, the Contradiction Engine might alert a monitoring service (through Interaction layer) that human oversight is needed. Or if the event was specifically a query from a user (sometimes the system might treat certain user queries as events too), it might directly proceed to answer using the updated memory. Generally, ingestion is decoupled from immediate response except for acknowledgement.

Throughout this flow, the guiding principle is maintaining a **low-entropy, high-coherence knowledge state**. By the end, the knowledge base either integrates the new information (lowering overall uncertainty by adding clarity) or discards it as noise, thereby keeping the entropy in check.

Query/Recall Flow (Retrieval and Response)

When a query is made to Kimera SWM (through the Interaction layer), the system uses a combination of direct lookup and associative activation to retrieve relevant information. The steps are:

1. **Query Interpretation:** The query arrives via the Interaction layer and is converted to an internal format (using EchoForm if needed). This internal query might be a structured template (e.g., subject=Skyship, relation=detected, object=Storm, find events) or a logical question (e.g., find contradictions in domain X, or summarize recent events about Y).
2. **Memory Activation (Hierarchical “Glowing” Patterns):** Once inside the Core, the query (especially if it’s a semantic or conceptual query) triggers a **spreading activation** in the memory structures. The Semantic Processor identifies key concepts in the query and “lights them up” in the Memory Vault’s graph. For instance, a query about “storm” will activate the concept node for “Storm” in the semantic network, giving it an initial weight. Through **hierarchical links** and associations, this activation spreads out to related nodes: e.g., “weather”, “rain”, or any event nodes linked to storms. This mechanism is akin to cognitive spreading activation, where source nodes are labeled with activation that propagates to linked nodes with decreasing strength ⁵. In Kimera, this results in a set of candidate memory entries (vault entries or scars) being “lit” or marked as relevant – these form a **glowing pattern** in memory, highlighting clusters of related knowledge.
3. **Focused Retrieval:** The system then focuses on the most strongly activated items. Hierarchically, if “Storm” as a concept is under a broader category “Weather Event”, that category might also glow and bring in any memory under it. The Entropy Manager might dynamically adjust this by boosting novel connections (to maximize information) or damping overly general ones (to avoid trivial info). The Semantic Processor queries the Memory Vault for entries that match the query constraints (e.g., all vault entries where subject or object = “Storm” and relation = “detected by Skyship”). Simultaneously or sequentially, it may query the Scars Repository for any recent events that fit, especially if the query asks for recent or specific instances.
4. **Aggregation of Results:** The system collects the results found. This could be an exact answer (if the query was specific and the vault had a matching entry, e.g., “Yes, Skyship detected a storm on 2025-06-01”) or a set of items that need synthesizing (if the query was broad like “What events happened today?” it might retrieve numerous scars from today). The glowing memory activation helps ensure that even indirectly relevant info is not missed – for example, if the query is “What could cause a storm?”, the concept “storm” might activate related causes like “weather system” or past events of storms along with their causes, giving the reasoning engine material to form an answer.
5. **Composition of Answer:** The Core now uses the retrieved information to formulate a response. If only one or few direct facts are needed, it will format those. If synthesis or reasoning is needed (like explanation or summarization), the core might perform additional processing: e.g., constructing a chain of reasoning from the causal links (using the mirror causal graph to trace backwards from an event to potential causes). Because the system stores cause-effect links, it can, for example, trace “storm detected” -> cause link -> “cold front collision” from a prior knowledge, thus answering a why-question. This composition is done in the Semantic Processor with awareness of the query type (factoid vs. explanatory vs. list).
6. **Output through Interaction:** The final structured answer is sent back to the Interaction layer’s Response Formatter, which then delivers it to the user in the requested format.

Non-Recall Mechanism in Retrieval: It is important to note that Kimera SWM does not treat all stored data as equally retrievable; it employs a **non-recall mechanism** to avoid surfacing irrelevant or low-confidence

information. This is akin to human memory inhibition – suppressing irrelevant memories so that relevant ones can be efficiently recalled ⁶ . In practice: - Some scars that were kept only for logging and flagged as low-entropy may be excluded from normal query results. The vault might only index scars of a certain significance for query answering. So even though they exist in storage, the system “chooses not to remember” them unless explicitly asked (e.g., a debugging mode or a specific query for full logs). - If contradictory information exists, one side of the contradiction might be suppressed in answers. For example, if the vault has two conflicting entries (one may be marked outdated by the contradiction engine), the query system will by default recall the *current* accepted one and not the refuted one. Only if specifically queried (or if the confidence of both is equal) would it mention uncertainty. This ensures the system’s responses remain coherent and not self-conflicting in normal operation. - Irrelevant associations that got activated in step 2 are dampened by the time of answer composition. The activation network naturally decays with distance ⁷ , meaning far-fetched memories won’t stay “glowing” enough to influence the answer. This is a design to keep output focused and is directly configurable via parameters like activation decay factor or a threshold to ignore weak activations.

Through these retrieval steps, Kimera SWM leverages both direct querying of structured knowledge and an associative recall process guided by activation and entropy. The combination allows precise answers when knowledge is directly stored and creative surfacing of related information when an answer needs more context or inference. All the while, entropy principles ensure that the recall doesn’t dredge up excessive or irrelevant data, aligning with the idea that **forgetting (or non-recall) is as important as remembering for efficient cognition** ⁶ .

Example Scenario

For illustration, suppose Kimera SWM has stored knowledge from previous events that “Skyship Alpha detected a Category 2 storm on 2025-06-01 near Location X” and “Skyship Beta reported clear skies on 2025-06-01 near Location X”. Now a user asks: **“Did the skyship encounter a storm?”**

- The query is parsed to an internal form: Query(subject=Skyship, verb=encounter, object=Storm?). The terms “Skyship” and “Storm” trigger activation in memory. “Skyship” may activate both Skyship Alpha and Beta entries, and “Storm” activates the storm event. The semantic network links “detected” ~ “encounter” as related concepts, so the storm detection event gets a strong activation as a candidate answer.
- The contradiction engine is not directly involved in this query, except that if the knowledge base had contradictory reports (storm vs clear sky), those entries might be marked with context (perhaps they are not directly contradictory since they could refer to different ships or times, but if they were same context, one might be flagged).
- The Memory Vault retrieval finds the fact of Skyship Alpha detecting a storm. It also finds Skyship Beta’s report (clear skies) because “Skyship” was a keyword, but that’s not directly relevant to “storm encounter” – however, since “storm” concept wasn’t in Beta’s report, its activation is weaker.
- The system composes the answer focusing on the storm event: “Yes, on June 1, 2025, Skyship Alpha encountered (detected) a Category 2 storm near Location X.” It might omit the Beta clear sky fact entirely due to low relevance (non-recall of irrelevant info) or it might include a nuance if asked specifically about that date and location, like “(Another skyship in the area reported clear skies, indicating the storm was localized.)” – that would be if the query was more detailed and the system decided the additional info is useful.
- The answer is returned, properly formatted in natural language.

This scenario shows how Kimera picks the right memory and suppresses others, and how stored details are used to form a coherent answer.

Advanced Memory Dynamics and Behaviors

Non-Recall and Memory Inhibition

Non-recall in Kimera SWM refers to the intentional non-retrieval of certain stored information unless specifically needed, mirroring the human ability to forget or ignore irrelevant details. This is a deliberate design to prevent information overload and maintain **focus** during reasoning. Technically, non-recall is achieved through a few mechanisms: - **Relevance Filtering:** Every query or cognitive task in the core is accompanied by a filter that excludes low-relevance items. For example, scars with an entropy score below a threshold are not even considered during retrieval, unless the query explicitly asks for exhaustive detail. They exist in storage for completeness but are effectively invisible to normal operations. - **Contextual Gates:** The system can mark certain information as context-specific. If outside the context, that info won't be recalled. For instance, facts relevant only to a certain scenario or time might be kept dormant when the context shifts. This is implemented by tagging vault entries with context labels and requiring a context match to activate them. - **Active Suppression of Confusing Data:** If a particular memory piece tends to produce confusion or contradiction, the system may suppress its recall proactively. This could be triggered by the contradiction engine – e.g., if two facts conflict, one might be labeled “inactive” pending resolution, so it won't be recalled in queries. - **Memory Inhibition Parameter:** At a higher level, an “inhibition” parameter could globally tune how aggressive the system is in ignoring less relevant memories. High inhibition means very selective recall (only the top few matching items), whereas low inhibition means broader recall. By default, Kimera is set to moderate inhibition to simulate efficient human-like recall (important things surface quickly, unimportant ones stay buried). As cognitive science suggests, such selective forgetting is adaptive for efficient recollection ⁶.

From an engineering perspective, non-recall is implemented in search and indexing algorithms. For example, the vault's search function might sort results by a relevance score (combining semantic similarity, recency, and entropy) and cut off after a certain score or number of results. The scars search might only look at the last N events or those above a significance level. These ensure that the system doesn't waste time on millions of trivial memories for each query. Tests can verify that as memory grows, query performance remains high thanks to this filtering (which is crucial for scalability).

It's important to note that non-recall does not mean deletion; the info is there and could be accessed if needed (for example, in a deep analysis mode or by explicitly querying archives). It simply means the default behavior is to **not recall** the irrelevant, much like a person doesn't actively remember every detail of every day unless prompted.

Hierarchical “Glowing” Memory Activation

The concept of hierarchical glowing patterns refers to the way memory nodes in the knowledge hierarchy light up in response to stimuli (queries or new information), and how this activation spreads in a controlled manner. This is essentially the **spreading activation** model applied to our semantic storage ⁸: - **Hierarchy in the Vault:** Knowledge in the Memory Vault is structured, often in a hierarchy or network. For example, “storm” is under “weather event” which is under “natural phenomena”. Also, “storm” might be linked to specific instances (like Storm event of 2025-06-01) and related concepts (like “rain” or “hurricane” if

Category 2 implies tropical storm). This creates layers of nodes from general to specific. - **Glowing Activation:** When one node is activated (say “storm” because of a query or because a new storm event was just stored), it glows – meaning we assign it a high activation value. Directly connected nodes (attributes of that storm, its category, the skyship that detected it) receive a slightly lower activation (the glow spreads). More distant nodes get a weaker activation and so on, often decaying exponentially with distance in the graph ⁵. The hierarchical structure ensures that not only the exact item is activated, but also its general category (which might tie in related memories) and its specific instances. - **Usage of Activation:** This glow has two main uses: 1. **Retrieval Augmentation:** As described in the query flow, activation helps retrieve indirectly relevant info. It’s essentially a smart search broadening: rather than only looking for exact matches to a query, the system looks at what’s glowing. For instance, if “storm” glows, it might retrieve not just the exact storm asked about, but other storms around that time or storms that affected the same location, if those nodes got some activation through shared links. This helps in providing context or detecting patterns (maybe the user’s implicit intent was to know if storms are common – multiple activations of storms would indicate that). 2. **Continuous Learning:** Activation can also guide learning. If a new event comes in and activates certain nodes, the system can strengthen the connections between those nodes. For example, each time “skyship” and “storm” are activated together (skyship detects storm), the link between the concept of skyships and storms is reinforced in the vault’s association network. Over time, this might form a pattern that the system recognizes (like “skyships often encounter storms”). So the glow patterns contribute to a kind of associative learning. - **Control of Spread:** The Entropy Manager plays a role here too. High entropy nodes (very novel or uncertain links) might be prevented from spreading activation too widely (since they could be noise). Low entropy, well-established nodes spread activation readily. Additionally, a **decay factor** is applied to prevent runaway activation: after the query is answered or after some time, the glow fades. This is analogous to short-term memory or working memory in humans – recently activated concepts remain briefly excited then return to baseline. It prevents the system from staying biased by an old query when a new query comes (avoid interference).

In implementation terms, this glowing activation could be done by maintaining an activation map in memory: a dictionary of node -> activation level that resets or decays over time or queries. Functions that search the vault combine keyword matching with activation values (e.g., score of a node = relevance to query + current activation level). The hierarchical nature means if a parent node is activated, all children might get a base activation too. This requires that hierarchical relationships are known; thus the vault likely has “is-a” or category edges that the activation algorithm follows.

Engineers can tune the glow parameters (initial activation strength, decay rate, etc.) and test that the system returns better results with it than without (e.g., recalling useful related info vs. missing them). It’s a flexible mechanism to emulate intuitive leaps or reminders that occur in human recall.

Mirror-Based Causal Propagation Logic

The mirror-based causal propagation logic is the framework ensuring that cause-effect relationships in the system are kept **consistent, bidirectional, and up-to-date**. The term “mirror-based” implies that for every causal link made in one direction, a mirrored link is maintained in the opposite direction, providing a two-way navigability in the event graph.

Key aspects of this logic include:

- **Bidirectional Linking:** Whenever the system records that Event A \rightarrow causes \rightarrow Event B, it automatically also records B \rightarrow is caused by \rightarrow A. This seems straightforward, but it's crucial. It allows queries like "why did B happen?" to directly find A, and queries like "what were the effects of A?" to find B. Without the mirror link, one of those directions would require a search through all events. By maintaining the explicit inverse relationship, retrieval of causal chains is efficient. The data structure for an EventScar might thus have two lists: `causes` and `effects`, each containing event IDs (populated in mirrored fashion).
- **Causal Graph Integrity:** The propagation logic also ensures there are no broken links or orphaned references. If an event is removed or archived, the system will propagate that change: any other event that listed it as a cause or effect is updated (e.g., remove the reference or mark it as an unknown cause). This is akin to referential integrity in databases. If a cause is retracted due to contradiction handling, any effect that depended on it might be flagged for review (since its explanation is gone).
- **Transitive Reasoning:** With mirror links, the system can do transitive traversal: if A caused B and B caused C, one can infer A indirectly influenced C. The propagation logic might choose to materialize some of these transitive links for quick access (for example, store A \rightarrow C as an "upstream cause" link). However, this could be exponential, so more likely it computes such chains on the fly when needed. Still, having the immediate mirror links simplifies building the chain.
- **Trigger Propagation:** Beyond just storing relationships, the causal propagation logic might allow events to trigger follow-up processing. For instance, if Event A is stored and it's known to cause type B events generally, the system could anticipate an Event B or at least be on the lookout. As a simple case, if input says "User clicked start, which causes system to initialize", the system, upon logging that cause, might proactively execute some initialization or expect an initialization event. While this edges into an active agent role, it could be part of how the system learns procedural or sequential knowledge.
- **Mirroring in Knowledge Representation:** The mirror principle can also apply to static knowledge in the vault: many relationships are inherently bidirectional (if X is part of Y, Y has part X; if $X > Y$ then $Y < X$ mathematically). Ensuring that whenever a relation is added, its logical inverse is considered or recorded if meaningful, helps the contradiction engine and query system. The vault's schema may define certain relations with automatic inverses (the EchoForm grammar could even specify this). For causal relations, it's explicitly done as above.

Implementation: The core likely has a Causal Propagation module (as described in Core layer) that is invoked whenever events are added/removed. It might use a publish-subscribe internally: e.g., the Scars Repository publishes "new event stored" with details, and the Causal module subscribes to update cause/effect links. The cause detection might require some logic or pattern matching – e.g., if the EchoForm content of an event has phrases like "because" or causal verbs, or domain rules (in a physical system, some events naturally cause others). Some causes might be inferred from temporal sequence or known dependencies (the knowledge base might know "if condition X then outcome Y" so if an event of type X happened, it infers a cause for a later Y event). These inferences are then stored as hypothesis links, possibly with lower confidence.

Use in Reasoning: When explaining or diagnosing, these causal links are invaluable. The system can trace backward from a problem event through mirrored links to find candidate causes. This is essentially automated causal reasoning. For consistency, the system may also propagate uncertainties: if a cause is

uncertain, it can mark the effect's explanation as uncertain. If a cause is contradicted (found to be false), the effect might inherit a contradiction tag (like "one of its recorded causes is false, so this event's explanation is in doubt"). Such propagation of flags helps ensure the knowledge doesn't become inconsistent unknowingly.

In summary, mirror-based propagation is about maintaining a **robust causal network**. It ensures that the knowledge graph of events remains symmetrical and navigable, enabling straightforward queries of "how?" and "why?" and underpinning the system's capability to maintain coherence even as events get updated or invalidated.

EchoForm Grammar as the Semantic Substrate

EchoForm is the custom grammar and representation language that underlies all of Kimera SWM's knowledge encoding. It serves as the **orthographic and semantic substrate**, meaning it defines both how things are symbolically written (orthography) and what they mean (semantics) in the system. This section details its role and structure:

Design of EchoForm:

EchoForm is essentially a formal grammar tailored to represent knowledge in a way that's both human-readable (to some extent) and machine-interpretable. It borrows concepts from structured languages (like logic or programming languages) but is aimed at capturing real-world events and facts. Key characteristics might include:

- **Explicit Structure:** An EchoForm expression has a clear structure, often resembling a logical assertion or a syntax tree. For example, an event could be represented as `(<EventType> [<role1>: <entity1>][<role2>: <entity2>]...` as we illustrated `("(DETECT)[subject: skyship][object: storm][time: 1500hrs]")`. This syntax is unambiguous to parse.
- **Orthographic Symbols:** It uses specific markers like parentheses, brackets, and keywords to denote roles and relationships. This is the orthographic aspect – a standardized writing system for knowledge. Because of this, EchoForm can capture things that normal NL sentences might obscure (like you can clearly mark the subject, object, causal link, etc., rather than relying on word order or prepositions that are language-specific).
- **Semantic Mappings:** Each EchoForm construct maps to a semantic concept in the system. For example, an event in parentheses indicates an instance of that event type in the ontology; roles like "subject" or "object" map to specific relationship types in the vault. In a sense, EchoForm could be seen as a domain-specific language for encoding a knowledge graph.

Role in Input Parsing: The Input layer's parser is essentially an interpreter or translator for EchoForm. If an input comes in natural language, the parser's job is to convert it to EchoForm. This might involve natural language processing or a simpler template matching if inputs are expected in a certain format. Because EchoForm is formal, the end result of parsing is something the core can rely on without ambiguity. If inputs are already structured (say JSON from a sensor), the adapter might directly produce EchoForm (like mapping JSON fields to an EchoForm statement). In some cases, users could input directly in EchoForm if they are technical (like writing a logic statement) – the system could accept that too.

Role in Storage: The Memory Vault might actually store knowledge in an EchoForm-like notation or at least use it as a serialization format. For example, a vault entry "SkyshipAlpha -> detected -> StormBeta" could be stored as `(DETECT)[subject: SkyshipAlpha][object: StormBeta][time: 2025-06-01]`. Alternatively, the vault might store a more normalized triple/graph form but can output EchoForm on demand. EchoForm thus acts as a lingua franca within the system – whenever knowledge is exchanged

between modules, it can be expressed in this form. This avoids misinterpretation and simplifies the interfaces.

Role in Output Generation: The EchoForm grammar can be reversed to generate natural language or other human-friendly output. The Response Formatter might take an internal EchoForm statement and have a set of templates or rules to read it off in English (or another language if extended). For example, `(DETECT)[subject: SkyshipAlpha][object: Storm][time: 3PM]` could map to a sentence "Skyship Alpha detected a storm at 3 PM." Because EchoForm captures roles explicitly, generating a sentence is straightforward by plugging in a sentence structure (subject-verb-object-time in this case). If multiple languages or formats are needed, one could create multiple renderers that consume EchoForm structures.

Orthographic Consistency: Using EchoForm means that the core logic doesn't have to handle multiple ways of saying the same thing. It reduces semantic entropy at the representation level – there's essentially one canonical form for a given piece of knowledge. This prevents issues like duplicate entries (e.g., "Skyship detected storm" vs "Storm was detected by skyship" would be normalized to one form). It simplifies contradiction detection as well: the engine can compare EchoForm strings or structured forms rather than doing complex NL understanding. It's easier to see that `(DETECT)[subject: A][object: B]` contradicts `(NOT-DETECTED)[subject: A][object: B]` if such a form is used for negation, for instance.

Extensibility of EchoForm: As the system's domain grows, EchoForm can be extended with new event types, roles, or ontological elements. It's akin to adding new words and rules to a language. The grammar is defined in a way that new productions (like a new event type "ALERT" or a new role "[cause: X]") can be added without breaking existing structure. The parser and generator need updates for new vocabulary, but the architecture remains the same. This is much easier than altering a whole NLP pipeline for a new domain, because EchoForm is controlled and formal.

Example: Suppose we want to add a concept of "location" to events (if not already present). In EchoForm, we might introduce a `[location: ...]` role. The grammar is updated: Event -> '(' EVENT_TYPE ')' ROLE_LIST and ROLE -> '[location: ENTITY]' becomes a possible role. The parser will then recognize "[location: Paris]" and map it to an attribute on the event object. The vault might then store location as an attribute of the event node. The generator could then produce sentences including location ("... at Paris"). This addition doesn't affect unrelated parts of the system thanks to the layered design; it mostly concerns the Input parser and perhaps minor tweaks in output.

In summary, EchoForm is a backbone for semantic clarity. It ensures every piece of information is encoded in a normalized, structured way at the point of entry, which dramatically simplifies downstream processing for consistency, reasoning, and communication. It's a powerful substrate tying together the orthography (how knowledge is written) with semantics (what knowledge means), thereby reducing ambiguity and error throughout the system.

Interface Definitions and Module Boundaries

To maintain **strict modularity**, each layer and major component in Kimera SWM communicates through defined interfaces. This section enumerates the key interfaces and the boundaries they enforce, which also points to how the system can be extended or tested module by module.

Input Layer to Core Interface

Name: `IngestEvent` (Core ingestion interface)

Description: Called by the Input/Ingestion layer to pass a newly parsed event to the Core. Could be a direct method call or a message on an internal bus.

Input: `event: EventScar` – The structured event object produced by the parser. Contains all necessary data (no raw text).

Output: Typically none (could be void if synchronous) or an ack status. The Input layer doesn't expect a result from Core; it's a one-way handoff.

Boundary enforcement: The Input layer does not manipulate memory or do any logic beyond parsing. It merely hands off to Core. Conversely, the Core does not parse – it assumes the event is ready. If the event is somehow malformed (violating EchoForm schema), the Core may reject it, but that would indicate an Input layer parsing bug. This interface allows easy insertion of a mock Input during testing – e.g., feed the Core with a series of EventScars to simulate a stream of inputs.

Core to Storage Interface

Name: Memory Storage API (Vault & Scars interfaces)

Description: Used by Core to store and retrieve knowledge. It can be implemented as an in-memory library in the same process (for zero dependency), but logically it's an interface. We break it down into two parts for Vault and Scars:

- `Vault.put(entry)` / `Vault.update(entry)` – adds or updates a VaultEntry. Might handle linking the entry in the graph, merging with existing if keys match, etc. Returns an ID or success status.
- `Vault.get(query)` – retrieves entries matching a query. The query might be a structured object or simply a key or concept. Returns one or a list of VaultEntry objects or a custom result object if we also want reasoning (e.g., could return a subgraph that answers the query).
- `Vault.remove(entry_id)` – (used by pruning or contradiction resolution) remove or deactivate an entry.
- `Scars.add(eventScar)` – appends a new event scar to the repository (like logging it). Returns the assigned ID.
- `Scars.find(criteria)` – search scars by filters (time range, type, related concept). Used for queries or propagation logic.
- Possibly `Scars.remove(id)` if we ever hard-delete, or `Scars.archive(id)`.

Boundary enforcement: Core logic never bypasses these calls – it doesn't, for example, write directly to a database or file. This means the storage mechanism can be changed (e.g., from in-memory to a database) by implementing these functions differently, without affecting core logic code. Also, it simplifies testing: one can replace the Storage API with a dummy in-memory stub that just records calls, to test core logic in isolation (or vice versa, test storage with simulated core calls).

Data Format: The interface uses internal data structures (EventScar, VaultEntry as described earlier) so both sides speak in terms of those. There's no JSON or external format at this boundary; it's all in-memory objects or structures, reinforcing zero external dependencies.

Core to Interaction Interface

Name: Query/Response API

Description: Functions or endpoints that the Interaction layer uses to send questions and commands to the core and receive answers. For example:

- `Core.answerQuery(queryObject) -> answerObject`
- `Core.getStatus()` - maybe to check health or summary of memory (for an admin interface)
- `Core.executeCommand(commandObject) -> result` - if the system allows imperative commands via interface (e.g., a command to purge something or load a dataset, though those might be internal or admin-only).

The primary one is the query answering. The `queryObject` would be an EchoForm structured query or some structured representation the core can use (this object could be produced by the Interaction's parsing of the user question). The `answerObject` might be an EchoForm structure for the answer or a simple string/number if the query was factual.

Boundary enforcement: The Interaction layer does not poke around inside core's memory or algorithms. It just uses these calls. That means the core can operate without any UI attached (for example, in a batch processing mode or during automated testing, one could directly call `answerQuery`). It also means multiple interaction layers could coexist (like a GUI and a CLI) both using the same core calls.

Stateless vs Stateful: Ideally, each query call is largely self-contained (the state is in the memory, not in a session). However, if we allow multi-turn conversational context, the core might accept a context handle or use some part of the vault as context memory. The interface could then have an identifier for session. For now, we consider it stateless to keep it simple – context is just part of the query object if needed.

Output Format: The `answerObject` could be raw data that the Interaction layer then formats. Or core could directly return a human-readable sentence. We have flexibility here, but in line with separation of concerns, it's cleaner that core returns structured info and let the Interaction format it. E.g., core might return something like `{ "result": True, "detail": "(DETECTED)[subject: SkyshipAlpha][object: Storm][time: ...]" }` and the Interaction layer turns that into "Yes, ...".

Internal Module Interfaces

Within the core, modules like the Contradiction Engine, Entropy Manager, etc., also interact through clear interfaces, even if they are just function calls: - The Semantic Processor might call `ContradictionEngine.check(newInfo)` which returns either `ConsistencyOK` or `ContradictionFound` along with details. - It might call `EntropyManager.assess(newInfo, context)` which returns an `entropy_score` and recommendation (like an enum: {KeepInVault, LogOnly, Reject}). - The Causal Propagation might be invoked as `CausalLinker.register(eventID)` which then internally accesses scars and vault to update links (the Semantic Processor doesn't need to manage that logic itself).

Each of these could be separate classes or services in the code, with their own internal data (the contradiction engine might keep a cache of known negations, the entropy manager might have some thresholds stored, etc.). They communicate by exchanging plain data (facts, scores, IDs) not by sharing global variables. This modularization means one can upgrade the Entropy Manager algorithm, for example, without altering how the rest of the core calls it – just keep the interface the same.

Extensibility Hooks

Because of the above interfaces, adding new capabilities often means extending an interface or plugging into one: - To **support a new input format**, one adds an Input Adapter for it, and perhaps extends the EchoForm grammar if new concepts are introduced. The Input->Core interface remains `ingestEvent`, so as long as the adapter produces a valid EventScar, the core doesn't care where it came from. - To **change storage engine**, one can replace the implementation behind `Vault.put/get` etc. If, say, moving from in-memory to a SQLite database for the vault, implement those functions to do DB operations. The core doesn't need changes aside from initializing the new storage module. - To **incorporate a new reasoning module** (maybe in the future, add a "Inference Engine" for deductive reasoning), we can have the Semantic Processor call it at a certain point. If it's optional, perhaps the interface is something like `InferenceEngine.tryDeduce(query)` that returns additional answers or nothing. If not present, the core can ignore it. This way, a plugin can be inserted without upheaval. - **APIs and Integration**: The Interaction layer can be expanded to provide new endpoints, e.g., a bulk data import function or an analytic query interface (like "export all events about X"). These would use combinations of core functions. Because core is self-contained, it could even be wrapped into different contexts (embedded in a larger application, or as a microservice).

The strong boundaries also help in **testing**: one can simulate each part. For instance: - Test Input parsing by feeding raw inputs and verifying the EventScar output. - Test Contradiction logic by giving it various existing vault states and a new fact, see if it flags correctly. - Test Entropy pruning by populating storage with many entries and calling a prune routine, check that high-entropy ones remain and low entropy removed. - Test full integration by using the Query API as a black-box: load some events, then ask queries, verify the answers.

By structuring the system in this decoupled way, Kimera SWM achieves a high degree of adaptability and reliability. Each module can be upgraded or replaced with minimal impact on the others, and the clear contracts (interfaces) between them mean miscommunications are minimized (a module can trust the format of data it receives). This is crucial in a complex system where new knowledge types and processing strategies might need to be introduced over time.

Entropy-Governed Memory Pruning Mechanism

One of the innovative aspects of Kimera SWM is its **entropy-governed pruning** – the way it manages memory growth and relevance by removing or archiving information based on entropy measures. We have touched on this in describing the Entropy Manager; here we consolidate how pruning works, the constraints around it, and its effects on the system.

Rationale for Pruning

Unlike traditional systems that might prune based purely on recency (e.g., keep the last N events) or frequency (remove least accessed items), Kimera SWM uses an information-theoretic approach: - Every piece of memory has an associated **semantic entropy** value (or inversely, an information value). This value can be thought of as how much uncertainty it reduces in the knowledge base or how much uniqueness it carries. - Over time, as more events accumulate, some information becomes redundant or less useful. If ten identical events occur, the latter ones add almost no new information beyond the first. Similarly, if a detailed event is fully distilled into the vault (all its useful info captured), the raw event record may become less necessary. - The system prunes to save space and to reduce “cognitive load” – a smaller active memory makes queries faster and contradiction checks easier (less clutter to sift through). However, it must do so without losing critical knowledge, hence the entropy-guided strategy to cut fat but not muscle.

How Pruning Operates

Pruning runs either periodically or triggered by certain conditions (like memory usage thresholds or when an entropy score is below a threshold on insertion). The process is roughly: 1. **Identify Candidates:** The Entropy Manager scans the memory for low-value entries. This could be: - Vault entries that have not been accessed in a long time *and* have low information content (e.g., they were highly expected or derivable from others). Possibly measured by an “information gain” metric when they were added. - Scars that have no links to any vault entries and have low novelty (e.g., an event that repeated known info and was not used in any reasoning). - Entries marked as superseded or contradicted (those are prime candidates since keeping them may cause confusion). The scan might produce a list of IDs with an entropy score and maybe a cost/benefit analysis for removal. 2. **Apply Thermodynamic Threshold:** The system could employ a rule akin to *Landauer’s principle metaphor*: avoid erasing info unless the entropy saved justifies it. In other words, only prune if it frees a significant amount of “memory entropy capacity”. This might translate to only pruning if the system’s total entropy is above a desired level or if the item’s entropy is below a certain fraction of the average. 3. **Prune or Compress:** For each candidate, the system decides one of: - **Archive and Remove:** The item’s essential info (if any) might be logged somewhere (external file or summary) then the item is removed from active memory. This is done for things truly not needed. Example: a raw sensor reading event that has since been aggregated into a daily summary – you drop the raw, keep the summary. - **Compression:** Sometimes instead of outright deletion, multiple low-entropy items can be merged into one. For example, if 100 identical events happened, the system might keep just one with a counter “100 occurrences” or convert them into a statistic. This retains the information (that it happened frequently) with far less storage. This is a form of entropy reduction by aggregation. - **Lower Resolution Storage:** Another approach is moving details to a slower storage. Since the core is self-contained, this might not be a separate DB, but conceptually, the system could mark an item as archived (not loaded in active memory) and store it in a compressed log on disk. It won’t be available for immediate recall (non-recall by design), but isn’t lost if needed offline. - **Retain (No action):** If for some reason an item is low entropy but policy says we keep a minimum record (like at least one event of each type for reference), the system may skip pruning it. 4. **Update Indices and Links:** After pruning, the relevant indices in Vault/Scars are updated. If a vault entry was removed, any references to it (like from scars or other entries) are cleaned. Often, the vault might not remove an entry entirely but mark it as *inactive* – effectively invisible to queries – to avoid re-indexing everything; then a background job could clean out inactive entries fully when convenient. 5. **Validation:** The system may run a quick consistency check (like the contradiction engine ensuring no dangling references or no loss of critical info). This is to confirm that pruning did not accidentally drop something important. The entropy calculation normally prevents that, but this check is a safety net.

Constraints and Safety

There are engineering constraints to ensure pruning doesn't harm the system: - **Never prune high entropy items:** By definition, those carry valuable info. Also, recent unique events shouldn't be dropped as they might be needed for upcoming queries. - **Avoid oscillation:** A danger could be if the system prunes something and later the same info comes, treating it as new again, oscillating. To avoid this, if something was pruned mostly because it was redundant, the fact that it was redundant means the knowledge is elsewhere in the system. If it wasn't actually redundant (e.g., we pruned a weird unique event by mistake), ideally we have it archived or the entropy threshold was mis-set. The system's design should be conservative – when in doubt, keep it (especially anything with potential future significance). - **User or Dev Control:** There should be some configuration for pruning aggressiveness. In a critical system, one might turn off automatic pruning and do it manually after review. Or set a policy like “don't prune events less than 1 day old” to ensure short-term memory is intact. - **Logging:** Any pruning action can itself be logged (maybe not into the same memory, that would be ironic, but into a maintenance log) so that developers can audit what was removed and why, which is helpful for debugging memory issues.

Impact on System Evolution

The entropy pruning mechanism allows Kimera SWM to scale and adapt over time without needing an ever-expanding memory. It tries to keep the memory **thermodynamically optimized** – maximum relevant information per unit of memory. This is especially important if the system runs continuously and could accumulate millions of events; with pruning, it can stabilize the knowledge base size or at least slow its growth to a manageable rate, focusing on meaningful data.

From a testing perspective, one would test the pruning by feeding many events and verifying that after prune, the remaining knowledge still answers queries correctly (no important regression). Also test that contradictory information, once pruned, indeed resolves the contradiction, etc.

In sum, entropy-governed pruning is a continuous housekeeping process in Kimera SWM that ensures the memory does not become cluttered with low-value information, thereby maintaining **high semantic fidelity** in the knowledge retained ⁹. It's a differentiator of Kimera's design, inspired by both information theory and cognitive efficiency, and is crucial for long-term robustness of the system.

Zero-Dependency Kernel Implementation

Kimera SWM's core is implemented as a **zero-dependency kernel**, meaning it avoids external libraries, frameworks, or services in its fundamental operations. This design decision has several implications on the architecture and engineering approach:

- **Language and Standard Library:** The system can be built in a performant systems language (like C++ or Rust, or even plain Python for a prototype) but without pulling in external packages. Data structures like graphs, trees, hash maps for the vault index, etc., are built using either standard library containers or custom implementations. This ensures that the behavior is fully under our control and there are no hidden processes (like an external DB engine) that could violate our entropy-based logic or cause unpredictable performance.
- **No ML Frameworks in Core:** For instance, if semantic similarity is needed (say to detect contradictions or to do query expansion), instead of relying on a heavy ML library at runtime, Kimera

might use precomputed embeddings or a simple vector math routine we implement. The rationale is to keep the core understandable and verifiable. If a neural network is needed for some semantic task, one might run it offline to generate data (like word embeddings) that the core can use in a lightweight manner. But the core itself won't depend on TensorFlow/PyTorch at runtime.

- **File Storage vs Database:** Instead of using a database server, the Memory Vault and Scars repository might be stored in a straightforward manner – possibly memory structures with periodic dumps to disk, or using a simple file-based index. This avoids complexity of external services and keeps latency low (no SQL query overhead, etc.). The trade-off is that we implement things like transactions or indices ourselves, but for moderate data sizes this is feasible.
- **Dependency Injection for Extensibility:** One place where we might allow an “external” component is via interfaces as discussed. For example, the Interaction layer could be considered outside the core and can use frameworks (like a web server library to expose an API). But the core doesn't rely on that – it could be tested with direct calls. The architecture allows the core to be a library that can be embedded anywhere.
- **Testing and Determinism:** With no external dependencies, the core's behavior is deterministic given the same input sequence (no randomness unless we introduce it, no external timing issues). This is great for testing reproducibility. It also simplifies reasoning about performance – we know exactly what algorithms are used for search, etc., so we can compute complexity and optimize if needed without being at the mercy of a black-box component.

Engineering Clarity: The zero-dependency approach enforces that every part of the system is explicitly defined in our codebase, which aligns with the requirement for strict engineering clarity. There's no “magic” from a framework that we don't understand – every function can be traced. This is especially important for something like the contradiction engine or entropy logic, where domain-specific custom behavior is needed that general libraries wouldn't provide.

Constraints: The constraint is that development effort is higher – we reinvent some wheels. But we consider it justified for a research/novel system like this where typical frameworks don't have the needed concepts (there's no off-the-shelf “semantic entropy memory” library). Moreover, keeping it lean on dependencies makes it easier to integrate into different environments (no complex install, easier to deploy on an embedded system if needed, etc.).

Kernel vs Extensions: We can think of the Core logic as the kernel. Around it, one could allow plug-ins or extensions that might use external tools. For example, perhaps an extension module for advanced natural language understanding could use an NLP library to parse text into EchoForm if our own parser isn't sufficient. That's okay because it's outside the core (in the Input layer). But once inside core, everything runs in our controlled code. Similarly, an extension for text-to-speech output would use external libs, but that's in the Interaction layer. This layered isolation means the critical memory functions remain isolated from those dependencies.

To ensure this in practice, a build of the core would not link against large frameworks. If using Python, we would not import heavy modules in core code (only possibly use small ones if absolutely needed, but aim to implement internally). Code reviews should trace all imports and references to confirm the dependency rule.

The benefit is a very **stable and transparent codebase** for the core. For long-term maintenance, this is a plus – fewer external updates to chase. For performance, often custom-tailored algorithms can outperform

generic ones for specific usage (e.g., our vault search might be super optimized for our query patterns, whereas a general RDF store might be slower).

In summary, the zero-dependency kernel means Kimera SWM's core is self-reliant and meticulously engineered from the ground up. This contributes to its clarity (every mechanism is explicitly defined), testability (no hidden interactions), and security (attack surface limited to our code). It's a conscious trade-off favoring control over convenience, aligning with the system's cutting-edge and experimental nature.

Extensibility and Adaptability for Evolution and Testing

The Kimera SWM architecture is built with future evolution and rigorous testing in mind. This section highlights how the design supports **extensibility** (adding new features or scaling up) and how components can be tested and tuned independently.

Extensibility Hooks

Throughout the design, we have pointed out places where new functionality can plug in: - **Adding New Memory Modules:** If we wanted to incorporate a new kind of memory (say an episodic memory module that stores raw transcripts of conversations, separate from semantic memory), we could do so by adding a new component in the Storage layer and interfacing it with the core. The core might be extended to consult this module for certain queries (maybe a query asks for exact recall, so it goes to episodic memory). Because our architecture already separates short-term event storage (scars) and structured memory (vault), adding another type is conceptually straightforward. - **New Analysis Modules:** Suppose in the future we want a **Temporal Reasoner** that can answer complex time-based queries or detect sequences. We could create such a module, feed it from the Scars repository, and have it communicate insights to the vault (like "Event X often follows Event Y"). This would integrate by having the Semantic Processor call it when needed or run it offline to annotate data. The key is, since core logic is modular, any new module just needs to adhere to the interface contracts. For instance, it could expose a method `findSequences(criteria)` that core can use when a query demands it. - **Domain Adaptation:** If Kimera SWM is deployed in a new domain (e.g., medical data vs. weather/robotics context), we can extend the EchoForm grammar with domain-specific vocabulary and add domain rules to the Contradiction Engine (like medical contradictions). The core algorithms remain the same; we just enrich the knowledge ontology. The modular design with EchoForm means the rest of the system will naturally support the new terms once they are parsed and stored. - **Scaling Up:** Should performance become a bottleneck, we have options to scale. One could distribute the Storage layer – e.g., have the Memory Vault on a separate process or machine (though that introduces a dependency). Because we have a well-defined interface (`Vault.get/put` etc.), that call could be turned into an RPC without changing how core uses it. Similarly, if the query volume is high, one could run multiple cores in parallel (since they're largely stateless between queries except shared storage). A load balancer in the Interaction layer could distribute queries among core instances that share a common storage backend. The architecture's separation of concerns allows for such scaling strategies. - **User Interface Changes:** As needs change, the Interaction layer can be adapted. If a visual interface is needed, it calls the same core API but perhaps queries more often or for different summaries. Or if integration with another AI is required (maybe Kimera SWM feeding a large language model), one can write an interface module that queries Kimera and feeds results to the LLM, or vice versa, without altering the memory core.

Adaptability for Testing and Tuning

The clarity of module boundaries makes the system highly testable: - **Unit Testing:** Each component can be unit tested with stubbed inputs. For example: - Feed the parser various sentences and ensure the EventScar output matches expected EchoForm structures. - Test the contradiction engine by constructing small vault scenarios (like a fact A, and then input not A) and verifying it flags it. - Test the entropy manager by giving it crafted events (with some known novelty values) and see if it recommends the right action. - Test the causal propagation by manually linking a chain of events and removing one, seeing if it cleans up. - **Integration Testing:** We can run the core end-to-end with a sequence of events and queries (simulating a real use case) and assert that the final state of the vault and the answers to queries match expectations. Because of determinism, these tests are reliable. - **Profiling and Tuning:** Since we control the algorithms, we can instrument them to measure performance (like how many vault entries are scanned per query, etc.). If something is slow, we can adjust data structures or add indexes. We're not dependent on an external system's black-box performance – everything can be traced. For instance, if queries slow down as vault grows, we might decide to add a caching layer for query results or optimize the spreading activation algorithm. These changes would be local to the core and transparent to the interfaces. - **Evolutionary Prototyping:** If an idea needs trying (say a new metric for entropy or a different strategy for memory activation), a developer can implement it in the corresponding module and see its effect. The system's modular design allows experimentation in one part without breaking others. We could A/B test two versions of the Entropy Manager by swapping it out via the interface (since the Semantic Processor just knows to call an `assess()` function; we can implement two versions of that easily). - **Logging and Monitoring:** We can add logging at module boundaries to observe system behavior in production. For example, log every contradiction detected, log pruning actions taken, etc. Since there's no external complexity, logs are straightforward to correlate with actions. This helps in validating that the entropy pruning, for example, isn't removing too much (we could run with logging enabled in a staging environment to see what it would prune, adjust thresholds accordingly).

Maintaining Clarity over Time

As the system evolves, maintaining the clear structure is important. One challenge in extensible systems is “feature creep” leading to tangled modules. To mitigate that: - We should continue to enforce that new features follow the layering. E.g., if we add a new reasoning capability, it should live in Core or be an adjunct service, but not sneak into the Interaction layer or bypass storage logic. - The EchoForm grammar should be the unified way to represent semantics. Even if we integrate some ML model or external knowledge, ideally we convert it into EchoForm/vault entries, rather than having parallel data forms. This keeps the substrate consistent. - Periodic refactoring might be needed as features add up, but thanks to interfaces, refactoring internals won't break the outer layers if done properly (e.g., we could change how vault data is structured internally – say switch from one graph library to another – and as long as `Vault.get/put` behave the same, everything else continues to work).

Example of Adaptation:

Imagine a future version of Kimera SWM wants to include a **probabilistic reasoning** module that can handle uncertain information (like Bayesian update of beliefs). We could implement this as follows: add a Probability Manager module in Core. It might hook into the Contradiction Engine (to handle degrees of belief instead of binary true/false) and into the Storage (to store probabilities with facts). Because our system was built assuming deterministic facts, we'd need to extend interfaces to pass probability info. This

could be done by extending VaultEntry with a `probability` field (the interface change) and updating ContradictionEngine.check to perhaps return a consistency measure. Thanks to modular design, these changes are localized – the Input layer wouldn't need changes (it can parse statements that include confidence perhaps), Interaction might format probabilities if present, but the overall architecture holds. The new module could reuse the event flow (maybe the Entropy Manager and Probability Manager work together: high entropy events get lower prior probability until confirmed, etc.).

We would test this addition by scenarios where contradictory info comes with different confidences and see that the system now doesn't outright discard one but labels probabilities appropriately.

In summary, the architecture's layered approach and clear interfaces make Kimera SWM **highly adaptable**. New features can be added with minimal disruption, scaling can be addressed by replacing modules behind interfaces, and the entire system can be systematically tested at multiple levels. This positions Kimera SWM as an evolving platform for semantic memory research and application, where improvements can be continuously integrated while preserving the core principles of entropy-guided, consistent, and efficient memory management.

1 Landauer's principle - Wikipedia

https://en.wikipedia.org/wiki/Landauer%27s_principle

2 Truth Maintenance System | by Rugvedi Ghule | Medium

<https://medium.com/@rugvedi.ghule20/truth-maintenance-system-28c7c2ef30f7>

3 4 9 Safari.pdf

<file:///file-3yz2knG9HiHXwAbvbYaQtX>

5 7 8 Spreading activation - Wikipedia

https://en.wikipedia.org/wiki/Spreading_activation

6 Memory inhibition - Wikipedia

https://en.wikipedia.org/wiki/Memory_inhibition