

# Todozi: AI-Driven Task Management System

## Whitepaper

### Introduction and System Overview

Todozi is an AI-first task management system that tightly integrates human workflows with AI assistance. It provides a **file-based** core (using JSON storage) for portability and simplicity, along with a rich set of features for project and task organization <sup>1</sup> <sup>2</sup>. The system is built to facilitate **AI/Human collaboration** at every level, from **structured task metadata** optimized for machine understanding to dedicated agent assignments and an **AI memory/idea system** for shared context <sup>3</sup>. Todozi can be used via a command-line interface (CLI) and an optional Terminal User Interface (TUI) locally, or as a networked service with a REST API – making it accessible to virtually any programming environment <sup>4</sup>. This whitepaper provides a technical deep-dive into Todozi's architecture and key components, aimed at AI researchers and developers interested in AI-augmented productivity tools.

At a high level, Todozi's design emphasizes **AI extensibility** and structured data. Tasks, ideas, and memories are stored in well-defined JSON schemas to be easily parsed and used by AI models. The system uses **semantic embeddings** to enable intelligent task search and recommendations beyond keyword matching <sup>5</sup>. A dedicated **embedding service** manages model loading, caching, and vector search operations. Todozi also incorporates a pipeline for **continuous model improvement**: collecting embedding data, exporting it for fine-tuning, and supporting model versioning and drift detection to ensure quality over time <sup>6</sup> <sup>7</sup>. Human users and AI agents collaborate through shared **context and memory** features – for example, AI can recall project-specific notes or prior decisions via a memory subsystem <sup>8</sup> <sup>9</sup>, and tasks can be extracted from natural conversation with the AI <sup>10</sup>. Finally, Todozi offers an interactive **Terminal UI** that presents tasks, analytics, and AI insights in a unified interface for efficient local operation. In the following sections, we explore each of these core areas in detail:

- **AI-Integrated Task Management** – how Todozi structures tasks, projects, and agents to enable AI involvement (metadata, assignments, hierarchical tasks).
- **Embedding and Semantic Search System** – the internal embedding infrastructure, including caching, semantic/hybrid search, clustering, recommendations, and drift detection.
- **Model Training and Enhancement Pipeline** – how embedding models are managed and improved over time (versioning, fine-tuning exports, validation).
- **Human-AI Collaboration Features** – memory and idea systems, shared context between humans and AI, and conversational task extraction.
- **Terminal User Interface (TUI)** – design and capabilities of Todozi's TUI for interactive task management and visualization.

Throughout the paper, code snippets and examples illustrate the implementation, and real-world use cases demonstrate how these features come together to support AI-augmented project management.

## AI-Integrated Task Management

Todozi's task management is designed from the ground up to accommodate both human users and AI agents working in tandem. At its core is a **structured task model** with rich metadata fields, enabling both fine-grained organization and machine readability. The **Task** structure (v1.2.0) includes fields such as unique ID, action (description), time estimate or deadline, priority, project association, status, assignee, tags, dependencies, context notes, progress, and timestamps <sup>11</sup> <sup>12</sup>. This comprehensive schema captures all relevant information about a task in a structured format. Notably, many fields are optional but standardized – for example, tasks can include an **assignee** indicating whether it's assigned to a human, an AI, a collaborative effort, or a specific AI agent role <sup>13</sup>. Tasks may also list **dependencies** (by referencing other task IDs) to express prerequisite relationships, and **tags** which serve for categorization and filtering <sup>14</sup>. This rich metadata not only helps human project managers but is also **optimized for AI processing**, as the structured JSON format can be easily serialized for model training or inference.

**Agent Assignment:** Todozi extends the concept of task ownership beyond just “human” or “AI” – it supports specialized **AI agents** that can be assigned to tasks. The system comes with default agent roles such as *Planner*, *Coder*, *Tester*, *Designer*, and *DevOps*, each intended for specific domains of work <sup>15</sup>. A task's assignee can be set as `agent:<role>` (e.g., `agent:planner`) to designate it to a specialized AI agent. Through the CLI, users can directly assign tasks to agents or switch assignees: e.g. `todozi assign <task_id> --assignee ai` to assign to the general AI, or `--assignee agent:planner` for a specific planner agent <sup>16</sup>. Under the hood, *Assignee* is an enum that includes `Ai`, `Human`, `Collaborative` (for tasks jointly handled by AI and human), as well as agent identifiers for specialized AI roles. This mechanism allows dividing work based on expertise – for instance, a “Coder” agent might handle programming tasks while a “Tester” agent focuses on QA tasks. The Todozi task list can be filtered by assignee to see which tasks are allocated to which agent or human (e.g., `todozi list --assignee agent:planner` to list all planning tasks) <sup>17</sup>. By tracking assignees in a structured way, Todozi enables clarity in AI-human responsibilities and supports downstream logic where agents can pick up tasks relevant to their role.

**Hierarchical Task Decomposition:** A standout feature of Todozi is how it supports breaking down complex projects into smaller tasks or “chunks” in a hierarchical fashion – a capability particularly useful when AI agents are involved in planning or coding tasks. Todozi implements a *chunking system* that allows tasks to be decomposed into multiple levels (project, module, class, method, block) with parent-child relationships <sup>18</sup> <sup>19</sup>. This was originally introduced to help an AI (e.g., a coding assistant) manage large coding projects by splitting them into bite-sized pieces. For example, a high-level project idea can be captured as a “project” chunk, which may be broken into several “module” chunks, each of which contains “class” or “method” chunks, and so on <sup>20</sup> <sup>21</sup>. Each chunk is identified by an ID and declares its level and an optional list of dependencies (other chunks it depends on) <sup>22</sup>. Todozi tracks these dependencies and can determine which chunks are “ready” (all prerequisites satisfied) so that AI agents focus on those first <sup>23</sup> <sup>24</sup>. In essence, this creates a dynamic task graph that the AI can navigate. The hierarchical decomposition enforces **token limits per level** (e.g., descriptions for project-level chunks are kept short, method-level chunks can have more detail) to ensure that each piece stays within manageable size for LLM context windows <sup>18</sup>. This structured breakdown vastly improves AI efficiency and scalability: the AI can maintain context at each level without being overwhelmed by the entire project at once <sup>25</sup> <sup>26</sup>. In real-world use, an AI *Planner* could create a project plan of chunks, and a *Coder* agent could then implement each code chunk sequentially, preserving context and respecting dependencies. Todozi's CLI and API support chunk management via commands like `todozi chunk list`, `todozi chunk graph` (to visualize dependency graph), and `todozi chunk ready` (to list ready-to-work chunks) <sup>27</sup>. This hierarchical task model not only mirrors

how human teams break projects into subtasks, but it also aligns with AI constraints, making Todozi particularly suitable for large projects involving AI co-workers.

**Rich Task Metadata for AI:** Every task in Todozi carries metadata that is leveraged to improve AI interactions. **Priority** and **Status** fields are enumerated types (e.g., Priority can be Low/Medium/High/Critical/Urgent; Status can be Todo, InProgress, Blocked, Review, Done, etc.)<sup>28</sup>. These structured fields allow an AI agent to prioritize actions (e.g., focusing on *urgent* tasks first) and to understand the state of each task in a predictable way. The optional **context\_notes** field provides additional natural language context or instructions for the task, which an AI agent can use to get more background or guidelines<sup>14</sup>. The **progress** field (0-100%) can be used to track partial completion of tasks, enabling AI to report how far along a task is (or even for an AI to update progress as it works)<sup>12</sup>. Because Todozi uses JSON for storage, an AI model fine-tuned on Todozi data can easily parse these fields. The **AI-first design** ensures that, for example, tasks added via natural language input are immediately structured in JSON with all these fields populated, so they can be fed into language models or embedding models for further reasoning<sup>3</sup>. In practice, a developer might say, “Implement user login; 3 days; high priority; project=website; status=todo; assignee=AI,” and Todozi would store this as a structured task that an AI agent can pick up. The structured metadata also makes it straightforward to perform bulk analysis – e.g., calculating how many tasks are assigned to AI vs humans, or how many high-priority tasks remain open.

**Project and Organization:** Todozi supports multiple projects, each essentially acting as a namespace or collection for tasks. Projects are stored in separate JSON files (e.g., `~/todozi/projects/myproject.json` for a project named “myproject”)<sup>29</sup>. This separation aids both humans and AI by providing contextual boundaries – an AI agent could filter its operations to a single project to avoid confusion with tasks from other projects. The CLI offers project management commands to create, list, archive, or delete projects<sup>30</sup>. Tasks themselves carry a `parent_project` field linking them to their project<sup>28</sup>. This ensures that when the AI or search system is asked about tasks, it can organize responses per project. Additionally, Todozi maintains global lists for active, completed, and archived tasks under the `tasks/` directory for quick access to tasks by status category<sup>31</sup>. The combination of project grouping and status grouping allows flexible querying (e.g., list all active tasks across projects, or list all tasks in project X that are blocked). AI agents can be aware of this structure; for instance, a *Planner* agent might first ensure a project exists, then populate it with tasks, and a *DevOps* agent might archive a project when it’s completed.

**Use Case Example:** Consider a software development team using Todozi. The team lead creates a new project “Website Launch” and outlines a few high-level tasks. Some tasks are assigned to human engineers, while others are assigned to AI agents. For example, a task “Set up CI/CD pipeline” might be assigned to `agent:devops`, while “Design landing page layout” is assigned to a human designer. The *Planner* AI agent, upon being assigned the project planning task, breaks it into chunks: a project-level chunk summarizing the whole website, module-level chunks for “Front-end”, “Back-end”, “Infrastructure”, etc., and further into class/method-level chunks for specific coding tasks. Each chunk appears as a Todozi task with a hierarchy indicated via dependencies (e.g., front-end tasks depend on the overall project chunk). As development proceeds, the *Coder* AI agent picks up method-level chunks (like “Implement user authentication method”) as they become ready, marking tasks InProgress or Done as appropriate. Meanwhile, the team can query Todozi for tasks by status or assignee – for instance, `todozi list --status blocked` to find tasks where the AI might be waiting on human input, or `todozi list --assignee agent:coder` to see what the coding agent is working on<sup>17</sup>. The structured nature of tasks means the AI agents can systematically coordinate: the *Tester* agent could list all tasks with status “done”

and automatically generate test cases for them, creating new tasks like “Test user authentication method” with itself as the assignee. This scenario illustrates how Todozi’s integrated task model facilitates a seamless division of labor, clear metadata for decision-making, and the ability to scale to complex, hierarchical workflows.

## Embedding and Semantic Search System

A core strength of Todozi is its **embedding-based semantic search** system, which imbues the platform with an understanding of content meaning. Instead of treating tasks and notes as plain text, Todozi converts textual content into high-dimensional vectors (embeddings) such that semantically similar items end up with nearby vectors. This enables features like finding related tasks even if they don’t share keywords, recommending next tasks based on context, and clustering tasks or ideas by topic automatically <sup>5</sup>.

**Embedding Model and Infrastructure:** By default, Todozi uses a pretrained transformer model (`sentence-transformers/all-MiniLM-L6-v2`) to generate 384-dimensional embeddings <sup>32</sup>. This model offers a good balance of speed and quality, and is downloaded and cached automatically on first use <sup>32</sup> <sup>33</sup>. The embedding generation process follows a typical BERT-based pipeline: input text is tokenized, fed through the transformer, then the token embeddings are mean-pooled and L2-normalized to produce the final dense vector <sup>34</sup> <sup>35</sup>. The resulting 384-float vector captures the semantic essence of the text. Todozi’s embedding subsystem is encapsulated in a service (`TodoziEmbeddingService`) that manages the model and a cache of embeddings in memory for fast reuse <sup>36</sup>. On initialization, the service loads the model (from the local HuggingFace cache under `~/.todozi/models/`) and prepares it for inference <sup>37</sup> <sup>35</sup>. The system ensures the model is only loaded once and stays in memory, resulting in quick subsequent embeddings (typical generation time ~50-100ms per task on CPU, with batching optimizations for multiple items) <sup>38</sup>.

One notable feature is the **embedding cache with smart LRU eviction**. Todozi keeps recently used embeddings in memory and can avoid re-computing vectors if the same text is encountered again <sup>39</sup>. The cache is bounded by a configurable memory limit, and least-recently-used items are evicted when space runs low <sup>40</sup>. This boosts performance in scenarios where certain tasks or queries repeat, achieving roughly a 60% cache hit rate in testing <sup>41</sup>. Moreover, embeddings are persisted in two places for durability: each task JSON includes its latest `embedding_vector`, and a global append-only log (`~/.todozi/embed/embedding_mega_log.jsonl`) stores every embedding with metadata <sup>42</sup> <sup>43</sup>. The latter serves as an “embedding diary” for analytics, backups, or model training data. An example log entry contains the timestamp, task ID, project name, original text, and the vector itself <sup>44</sup>. This design ensures that even if the model is updated or the cache is cleared, historical embeddings and their contexts are retained.

**Semantic Search and Hybrid Search:** With embeddings in place, Todozi enables powerful semantic search across all content types. A user (or agent) can search by meaning – for example, find tasks related to “login” even if the task description uses the term “authentication”. The CLI provides commands like `todozi search "<query>"` which under the hood will use semantic similarity to rank results <sup>45</sup>. There is also `todozi similar "<text>"` for finding tasks similar to a given piece of text, as demonstrated in the CLI example where searching for “add login system” finds a task about OAuth2 even though the keywords differ <sup>46</sup>. Internally, the embedding service computes the vector for the query and then computes cosine similarity with stored task vectors to retrieve the top matches. Cosine similarity is used as the metric (since

vectors are L2-normalized, cosine similarity reduces to dot product), and a score of 1.0 means identical semantic content while lower scores indicate less relevance <sup>47</sup> <sup>48</sup> . This allows Todozi to rank results by semantic relevance rather than lexical matching.

In addition, Todozi supports **hybrid search**, combining semantic similarity with traditional keyword filtering for even more precise results. A hybrid search takes into account both the embedding proximity and keyword overlap. For instance, an agent or user can perform a query that boosts results which contain certain tags or terms, while still using semantics to catch broader context. The embedding service offers a method `hybrid_search(query, keywords, weight, limit)` where a weight (0.0 to 1.0) is used to balance semantic score vs keyword score <sup>49</sup> . In practice, one might set 70% semantic and 30% keyword weighting, which was found to improve search accuracy by ~40% in tests <sup>50</sup> . The code snippet below illustrates a hybrid search example using the Rust API:

```
// Perform a hybrid search: "authentication" query, with a keyword filter
"security"
let results = service.hybrid_search(
    "authentication",
    vec!["security".to_string()], // keywords to match
    0.7,                          // 70% semantic weight, 30% keyword
    5                             // limit to top 5 results
).await?;
for r in results {
    println!("Found {} (score {:.2}) - {}",
        r.content_id, r.similarity_score, r.text_content);
}
```

In this example, the search will return tasks or notes related to authentication but prioritize those that also mention security, mixing the strengths of both approaches <sup>51</sup> <sup>52</sup> . This feature is particularly useful in large projects where you might remember a specific term associated with a task but still want the semantic breadth of embedding search.

**Clustering and Knowledge Graphs:** Beyond search, Todozi leverages embeddings to analyze and organize content globally. One feature is **automatic clustering** of related content. By analyzing distances between all pairs of content vectors, the system can form clusters of tasks, ideas, or code chunks that share thematic similarity <sup>5</sup> <sup>53</sup> . The embedding service provides a `cluster_content()` method which groups items and even computes summary statistics per cluster (such as average similarity within the cluster) <sup>53</sup> <sup>54</sup> . A cluster might represent, for example, all tasks and ideas related to “user authentication” or “database optimization.” These clusters can help both AI and humans identify hotspots in the project or areas of related work. Todozi’s TUI (if enabled) exposes visualizations for cluster overviews – for instance, showing how many clusters exist and their thematic labels or key terms (the system can auto-label clusters by common keywords or by nearest centroids, using types like `LabeledCluster`) <sup>55</sup> <sup>56</sup> .

Furthermore, Todozi can construct a **similarity graph** (a kind of lightweight knowledge graph) where nodes represent content (tasks, memories, etc.) and edges connect items with high similarity. The embedding enhancements introduced data structures like `SimilarityGraph`, `GraphNode`, and `GraphEdge` to represent these relationships <sup>57</sup> . The method `build_similarity_graph(threshold)` generates a

graph linking all items that have a cosine similarity above a certain threshold <sup>58</sup>. This is useful for uncovering implicit connections in the knowledge base – for example, if two tasks in different projects are similar, they would be connected, indicating a potential duplicate effort or a concept overlap that the team should be aware of. The TUI is capable of visualizing such graphs, allowing users to navigate the network of related items and possibly merge or cross-reference them.

**Recommendations and Next Task Prediction:** Another AI-powered feature is Todozi's **recommendation engine** for tasks and content. Using methods like `recommend_similar(based_on, exclude, limit)`, the system can suggest tasks or ideas that are related to a given set of items <sup>59</sup>. For example, if a user has completed tasks `A` and `B`, Todozi can suggest the next tasks that are semantically similar to `A` and `B` (perhaps indicating continuing work in that area) but not already done or in progress (using an exclude list) <sup>59</sup>. This works by calculating an "interest centroid" in the embedding space – essentially averaging the vectors of the *based\_on* tasks to find the general theme, and then retrieving new items nearest to that centroid <sup>60</sup>. The idea is to mimic how a human might think "Since you worked on X and Y, you might be interested in Z." Agents in Todozi can call `recommend_similar` to help planning; for instance, a *Planner* agent could use it to propose next steps for a user, or a *Coder* agent could find which unimplemented features are related to what's currently being built <sup>61</sup>. In testing, this approach provides intelligent suggestions, especially when the project has many loosely connected tasks – it helps surface items that otherwise might be forgotten. The recommendation system also takes into account content *types*, meaning it could suggest an idea or memory relevant to the task at hand (for example, reminding the user of a design idea that matches the current work). This cross-content awareness is a direct result of using a unified embedding space for tasks, ideas, memories, etc., so recommendations can span these categories.

**Drift Detection and Content Evolution:** Over the lifecycle of a project, tasks and content can evolve – descriptions may be updated, specifications may change. Todozi includes a **drift tracking** mechanism to monitor how embeddings change over time for a given item. Each time an item's text is significantly modified, a new embedding is generated, and the method `track_embedding_drift(id, new_text)` can compare the new vector with the original or prior vectors <sup>62</sup>. It computes a drift percentage (for example, based on cosine distance) and flags if the change exceeds a threshold (default >20% change is considered significant drift) <sup>63</sup>. The system can log these changes in a `DriftSnapshot` and compile a `DriftReport` that shows how a task's meaning has shifted <sup>57</sup>. This is important for long-running projects or living documents – if a task's embedding drifts a lot, it might indicate the task's scope has changed or it has been repurposed. In an AI context, this signals that the AI might need to re-evaluate related items; for example, if a "UI redesign" task drifted significantly after some discussion, any recommendations or clusters involving that task should be updated. Todozi's embedding service can maintain versioned embeddings for content: using `create_embedding_version(id, label)`, it snapshots the current embedding with a label (perhaps a date or change description) so that a history is kept <sup>7</sup>. Developers can retrieve this history via `get_version_history(content_id)` to see all past versions of an embedding <sup>64</sup>. These capabilities essentially bring *configuration management* to semantic knowledge – ensuring the system is aware of and can account for content drift. The TUI can display drift reports, and identify items that have changed meaning recently, which might require attention.

**Performance and Caching:** The semantic search system is built for efficiency. As noted, the model is loaded once and reused, with on-disk caching of model files to avoid re-downloads <sup>65</sup>. Batch operations are used whenever possible; for instance, Todozi's embedding service offers `generate_embeddings_batch(texts)` which can embed many items in parallel (using Rust's async tasks) for a ~10× throughput improvement over single-item embedding <sup>66</sup>. Search operations are

optimized by limiting vector comparisons to relevant subsets: if a user searches within a project, only that project's vectors are scanned, and even global searches can be accelerated by approximate methods or indexing (though specifics depend on implementation; the documentation notes average search latency <5 ms for typical usage with caching) <sup>41</sup>. The LRU cache not only avoids recomputation but can also cache search results for recent queries. For example, if an agent frequently queries related to "database", Todozi may keep a short-term index of top results for "database" so subsequent calls are near-instant. Memory usage is moderate: the model itself takes ~200 MB RAM loaded, each embedding vector ~1.5 KB, so even thousands of tasks consume only a few MB for their embeddings <sup>67</sup>. This footprint is acceptable for most developer machines or servers.

**Extensibility of Embedding System:** Todozi's embedding system is not a black box; it's designed to be extensible and transparent for developers: - **Custom Models:** Users can switch the embedding model at any time to better suit their domain or language needs. The CLI command `todozi emb set-model <model-name>` pulls a model from Hugging Face and sets it as the new default <sup>68</sup> <sup>69</sup>. For example, one could switch to a larger model like `all-mpnet-base-v2` for higher accuracy or a multilingual model for international projects <sup>70</sup>. Todozi will download and cache the new model and update its config (`tdz.hlx` config file) accordingly <sup>69</sup>. Multi-model support is also planned/available – the system can load additional models concurrently via `load_additional_model(name, alias)` and even compare them side-by-side <sup>71</sup>. This is useful for benchmarking a new model versus the existing one: using `compare_models(text, [model1, model2])`, Todozi can output the two embeddings and measure differences, including performance timings and quality metrics <sup>72</sup>. - **Analytics and Debugging:** To ensure the embeddings remain high-quality, Todozi includes a validation routine. The `validate_embeddings()` method scans through stored vectors to detect anomalies like NaNs, zero vectors, or outliers that deviate from the expected distribution <sup>73</sup>. It produces a `ValidationReport` highlighting any issues so that developers can address them (e.g., by regenerating certain embeddings or investigating input anomalies) <sup>74</sup>. There's also an `explain_search_result(query, result)` function which can help interpret *why* a given result was retrieved for a semantic query – it provides a breakdown of the similarity contributions and even identifies which vector dimensions contributed most to the match <sup>75</sup>. Such transparency is valuable for AI researchers to trust and fine-tune the system's semantic behaviors. - **Integration Points:** The embedding service is utilized across Todozi's subsystems: the CLI search commands and the REST API endpoints (e.g., `GET /embeddings/search`) call into it, the agents use it for finding relevant info (agents can call `hybrid_search` or `recommend_similar` as helpers), and the TUI uses it to display visualizations like similarity graphs and cluster maps <sup>76</sup> <sup>77</sup>. This consistent integration means improvements to the embedding engine benefit the entire system uniformly.

In summary, Todozi's embedding and search system transforms the way tasks and knowledge are handled by enabling semantic understanding. Instead of rigidly defined tasks, the system "understands" tasks in context – a crucial difference when scaling up AI assistance. Real-world use cases include a developer quickly finding all tasks related to a feature by description (even if phrased differently), an AI agent automatically grouping related tasks to suggest if one completion could close out others, or detecting when the project's focus has shifted because many tasks drifted semantically in a new direction. These capabilities greatly enhance productivity and insight in an AI-augmented project workflow.

## Model Training and Enhancement Pipeline

While Todozi starts with powerful pre-trained language models for embeddings, it also provides a framework to continuously **train, adapt, and improve** these models to better fit a team's unique data. This pipeline covers model version management, fine-tuning, and quality assurance to ensure the AI components remain effective as usage grows.

**Model Management and Versioning:** Todozi employs a versioned approach to managing embedding models. The current model in use is recorded in a config file (`~/.todozi/tdz.hlx`) under the `[embedding]` section, e.g., `model_name = "sentence-transformers/all-MiniLM-L6-v2"` <sup>78</sup>. When a new model is set via `todozi emb set-model`, this config updates, and the service knows to load that model on next initialization. Todozi also keeps models in a registry which supports aliases and multiple models loaded concurrently for comparison or A/B testing <sup>71</sup>. This is achieved with the `load_additional_model(name, alias)` function that loads another model (without switching the default) and stores it in memory with an alias for reference <sup>71</sup>. Developers can then use `compare_models(text, [alias1, alias2])` to generate embeddings from both models and compare results side by side <sup>72</sup>. The comparison includes timing (inference speed), output vector differences (e.g., cosine similarity between the two model outputs for the same text), and dimension consistency checks <sup>79</sup>. This facility is extremely useful when iterating on model improvements – one can quantify how a new fine-tuned model diverges from the old and whether it truly captures domain-specific nuances better.

In addition to model versions, the embeddings themselves are versioned per content. As mentioned, Todozi can snapshot embeddings over time using `create_embedding_version(id, label)` <sup>7</sup>. Each snapshot appends an entry (with the label and timestamp) to a version log (likely stored as JSONL similar to the mega log) for that content item. By retrieving the version history, one can visualize how an item's embedding moved in vector space over chronological edits <sup>64</sup>. This is especially relevant if the underlying model is changed – the next time an item's embedding is generated with a new model, that could be logged as a new version with a label like "after\_model\_upgrade\_v2". In essence, Todozi treats embeddings as data that can be migrated and tracked, which is important for long-term maintainability of AI features.

**Fine-Tuning Data Export:** As users accumulate tasks, ideas, and other content, they generate a valuable dataset that can be used to fine-tune embedding models to their specific domain or vocabulary. Todozi facilitates this by allowing the export of training data. The command or method `export_for_fine_tuning(output_path)` will collect all relevant text and metadata from the Todozi database and write it in a machine learning friendly format (JSONL) for model training <sup>80</sup> <sup>81</sup>. Specifically, it includes fields like the text content, current embedding vector, and possibly context such as tags or project (depending on how the fine-tuning is meant to be supervised) <sup>80</sup>. This output can then be used to fine-tune the sentence transformer model (or any embedding model) using frameworks like PyTorch or TensorFlow. By providing both the text and its embedding, Todozi can allow a fine-tuning process to use a form of knowledge distillation or contrastive learning: the current embeddings serve as a target or baseline. Alternatively, if tasks are labeled or categorized (e.g., by project or priority), those labels could be used in a supervised fine-tuning objective. The key point is that Todozi lowers the barrier to improving the model with one's own data by preparing the dataset in the right format <sup>81</sup>.



Once a model is fine-tuned externally (say, producing a new model file), integrating it back is straightforward: one would use `todozi emb set-model /path/to/new/model` (if it's a local model path) or register it in HuggingFace and refer to its name. Because Todozi uses the HuggingFace model hub mechanism, a fine-tuned model can be loaded either from a local directory or from a hub repository as long as it's compatible (i.e., same architecture). This means teams can iterate on models and deploy updates to Todozi without changing the codebase – just by swapping config.

**Continuous Improvement Loop:** The envisioned use of the above features is a continuous improvement loop for the AI models:

1. **Data Collection:** As Todozi is used, the embedding mega-log grows with diverse examples of tasks, ideas, and their embeddings. Users might also add explicit ratings or notes on task suggestions (though not described in the core docs, such feedback could be captured in future to supervise learning).
2. **Quality Monitoring:** Regularly, a developer can run `validate_embeddings` to ensure nothing is off in the current data <sup>73</sup>. Also, using `explain_search_result` on some queries can highlight if the model's understanding aligns with expectations (for example, if weird dimensions are dominating similarity, it might indicate a quirk to address) <sup>75</sup>.
3. **Fine-Tune Preparation:** After accumulating sufficient new data or observing areas where the model could do better (e.g., perhaps the model isn't distinguishing well between certain technical terms), the team exports the data via `export_for_fine_tuning` <sup>81</sup>. This yields a training file containing examples of the content and possibly their context.
4. **Fine-Tuning & Testing:** The team fine-tunes the model using standard training scripts, then tests the new model on some tasks – possibly using Todozi's `compare_models` to see differences in embedding outcomes for known queries or using `todozi similar` to subjectively see if results have improved.
5. **Deployment:** If the new model is deemed better, they load it into Todozi (`set-model`). At this point, Todozi's versioning can snapshot that a model change occurred. They might use `backup_embeddings` to save all current embeddings before the switch <sup>82</sup>, then re-embed content with the new model. Because all tasks can be re-embedded (perhaps via a re-index or on-the-fly as they are accessed), the system can smoothly migrate to using the new model's vectors everywhere. If needed, the old embeddings can be restored or compared if something goes wrong (using `restore_embeddings` if they saved a backup, or just by comparing drift on each item to see how much the new model changed things) <sup>83</sup>.
6. **Monitoring Post-Upgrade:** After switching models, one can again run drift detection globally to ensure no critical tasks lost their similar neighbors or that clusters didn't break apart unexpectedly. If the fine-tuning was successful, presumably similar items remain similar or even tighter clustered, whereas if the model introduced anomalies, the validation might catch an outlier embedding.

Todozi's commitment to **quality validation** is evident in the features like `ValidationReport` which can flag anomalies <sup>55</sup>, and the logging of performance metrics and diagnostics in structures such as `PerformanceMetrics` and `DiagnosticReport` <sup>84</sup>. For instance, performance metrics could record average embedding time, cache hits, etc., before and after a model update to quantify improvement or regression.

**Enhancements Summary:** The embedding subsystem of Todozi received significant enhancements (as referenced in an *Embedding Enhancements Summary*), adding 27 new methods and 11 new data structures <sup>85</sup> <sup>86</sup>. These were aimed at making it “enterprise-grade,” including batch processing, hybrid search, caching, cross-content search, clustering, recommendations, drift tracking, multi-model, and more. After these enhancements, the system boasts 50 public async methods for embeddings <sup>87</sup>. The development process included thorough testing (`cargo test --lib emb`) and an example demo (`embedding_enhancements_demo.rs`) to showcase everything <sup>88</sup> <sup>89</sup>. The result is a **production-ready** embedding engine with high-performance characteristics (10× batching speedup, <5 ms search latency) and

robust error handling <sup>41</sup>. For developers and researchers, this means Todozi's AI backbone is not a toy model but a scalable component that can be studied, extended, and integrated into larger AI systems. It could even serve as a blueprint for embedding services in other applications.

In summary, the **model training and enhancement pipeline** in Todozi ensures that the AI capabilities can evolve over time. Instead of being stuck with a static model, Todozi allows adaptation to the specific task domain and continuous learning. This is critical in real-world use: imagine a team of legal researchers using Todozi – they could fine-tune the embeddings on legal text so that semantic search understands that domain's vocabulary. Or a game development team might fine-tune on past game design documents to better cluster and search new ideas. Todozi provides the hooks to do this safely (with version control and backups), encouraging an *AI that grows with your project*.

## Human-AI Collaboration Features

Collaboration between human users and AI agents is at the heart of Todozi's design. The system includes dedicated features to share context, preserve knowledge, and facilitate natural interactions (like conversations) between humans and AI. These features ensure that the AI is not a black-box working in isolation, but rather a participant that can **remember, ideate, and converse** in alignment with the human team's goals.

**Shared Memory System:** Todozi introduces an AI **memory system** that allows storing contextual notes known as "memories." A memory in Todozi represents a piece of information that might be relevant for future tasks or decisions – essentially, a knowledge snippet that the AI can recall later. Each memory is a structured record with fields such as: moment (a short title or timestamp/context), meaning (the content or takeaway), reason (why it matters), importance (priority level), and term (short-term or long-term) <sup>90</sup>. For example, a memory could be: *Moment*: "2025-01-13 Client Meeting", *Meaning*: "Client prefers iterative development approach", *Reason*: "Impacts testing cycles and delivery schedule", *Importance*: high, *Term*: long-term <sup>91</sup> <sup>92</sup>. This might be something an AI PM or the human team lead wants to remember throughout the project. Todozi's CLI provides commands to create and manage these memories: `todozi memory create --moment "<...>" --meaning "<...>" --reason "<...>" --importance <...> --term <short|long>`, as well as specialized variants like `create-secret` (for AI-visible-only secrets) and `create-human` (explicitly marked as human-visible) <sup>93</sup> <sup>94</sup>. Memories are stored as individual JSON files in `~/.todozi/memories/` (named by UUID) <sup>95</sup> <sup>96</sup>. They are also indexed by the embedding service: when a memory is created, the text of its fields is embedded and added to the semantic search index <sup>97</sup>. This means that later on, an AI agent can semantically search through memories to recall relevant info. For instance, if an AI agent is planning a task related to testing, it might query memories for "testing approach" and retrieve the above memory about iterative development, which could influence its plan.

Memories have an associated **visibility and type**. The *standard* memories are accessible to both AI and human (e.g., shown in CLI or UI lists), whereas *secret* memories are intended only for AI (not listed in normal human views, perhaps containing sensitive data like an API key or an AI's internal note) <sup>98</sup>. Conversely, *human* memories might be things the AI should not use unless explicitly allowed (though the exact distinction in usage might depend on policies). There is also an *emotional* memory type that includes an emotion field, allowing the system to store a feeling or affective context (this is part of Todozi's extensibility to capture more nuanced context) <sup>99</sup>. For example, an emotional memory could record that "Team morale was low after missing a deadline" with an emotion "frustration" intensity 7/10. This could be used by AI

agents to modulate their suggestions or interactions (a very cutting-edge feature bridging into affective computing).

The memory system ensures **context persistence** across AI interactions. If a conversation with an AI agent results in a decision or an important fact, encapsulating it as a `<memory>...</memory>` entry means it gets saved and is retrievable later <sup>91</sup> <sup>92</sup>. The integration of memory in chat is seamless: as shown in the chat example, a user can say “I need to remember this: `<memory> ... </memory>`” and the `process_chat_message_extended` function will parse and save that memory entry <sup>100</sup> <sup>91</sup>. This is akin to giving the AI a notepad where both the user and AI can jot down things to remember. When making decisions or generating plans, the AI can be programmed to check the memory store for any relevant entries (Todozi facilitates this via functions like `list_memories` with filters, and semantic `search_memories(query)` that works like task search but on the memory database <sup>101</sup> <sup>102</sup>). Real-world use case: In a long project, the team might accumulate dozens of decisions and preferences as memories – the AI can then answer questions like “Why are we using PostgreSQL instead of MySQL?” by retrieving the memory where that decision was made, providing rationale that was stored (e.g., “decision made on X date for Y reason”). This avoids re-hashing discussions and helps new team members or AI agents quickly get up to speed on project context.

**Idea Management:** Alongside tasks and memories, Todozi has an **idea management** feature for capturing brainstorms, high-level concepts, or future to-dos that aren't yet tasks. Ideas are less structured than tasks (usually just a short text plus some metadata) and can later be turned into tasks or projects. The idea model includes fields like the idea text, an importance level, optional tags, optional context description, and a *share level* (private or public) <sup>103</sup> <sup>104</sup>. **Share level** determines visibility and who/what can act on the idea: a *private* idea is only visible to the user (and possibly their personal AI agents), whereas a *public* idea is shared more broadly, e.g., if the Todozi system is used by a team or connected to a central server, public ideas might be accessible via the API to other collaborators <sup>105</sup>. The CLI commands for ideas include `todozi idea create --idea "..." --share <level> --importance <...> --tags <...> --context <...>` to add an idea, and listing or showing them similar to tasks <sup>106</sup> <sup>107</sup>. Internally, ideas are stored as JSON in `~/.todozi/ideas/` with their UUIDs <sup>108</sup>. On creation, like tasks and memories, an idea's text is embedded for semantic search indexing <sup>109</sup>. This ensures that, for example, using the search feature can find relevant ideas when searching tasks (Todozi has a unified search or a `search-all` that can include ideas, tasks, memories, etc. in the results) <sup>110</sup>.

Ideas serve as a **collaboration bridge** in many ways. They allow a human user to input something that's not yet a formal task (“maybe use microservices for this component”) and later the AI can help flesh it out into tasks. Conversely, an AI agent might generate ideas (especially a creative or brainstorming agent) which the human can then review and promote to tasks or discard. The concept of “share level” also indicates a trust or publishing mechanism: for instance, an AI could mark an idea as `share: public` if it's confident and wants to propose it to the whole team, whereas it might keep an idea `private` if it's a rough thought it's unsure about. Todozi's server exposes ideas via endpoints (`GET /ideas`) so that a web UI or other clients can display shared ideas, facilitating team brainstorming in an AI-assisted manner <sup>105</sup> <sup>111</sup>. In a real scenario, imagine a team meeting where an AI listens and directly logs ideas as they come up: “<idea>Implement dark mode; share public; importance: medium; tags: UI, theme; context: users have requested this feature</idea>”. This would create an idea entry <sup>107</sup> <sup>104</sup> which is now searchable and can be turned into a task via a command or UI action. By capturing these in a structured way, nothing gets lost, and the AI can later remind the team of these ideas or even begin planning them when current tasks are completed.

**Conversational Task Extraction (Chat Integration):** One of Todozi's most innovative collaboration features is its ability to parse natural language conversations (chat) and **extract structured tasks, ideas, and other items** from it. This is achieved through a simple markup in messages: special tags like `<todozi>...</todozi>`, `<memory>...</memory>`, `<idea>...</idea>`, `<todozi_agent>...</todozi_agent>`, `<chunk>...</chunk>` encapsulate different types of content in a chat message <sup>100</sup> <sup>112</sup>. For example, a user or an AI could say: "Let me assign some tasks: `<todozi_agent>planner; task_001; project_planning</todozi_agent>` ... And create the tasks: `<todozi>Design microservices architecture; 2 weeks; high; system-design; todo; assignee=agent:planner; tags=architecture</todozi>`" <sup>113</sup> <sup>114</sup>. When this message is processed by `process_chat_message_extended`, Todozi will create an agent assignment (linking an agent "planner" to a placeholder task id in `project_planning`) and a new task "Design microservices architecture" with the given properties (2 weeks, high priority, etc.), assigned to planner agent, tagged architecture <sup>113</sup> <sup>114</sup>. Similarly, `<memory> ... </memory>` segments cause new memories to be saved, and `<chunk> ... </chunk>` lines produce code chunk entries in the chunking system <sup>91</sup> <sup>92</sup>. The output of processing a chat can be seen in the JSON response from the REST API: it lists how many tasks, memories, ideas, etc., were extracted and provides their details <sup>115</sup> <sup>116</sup>. For instance, after processing, it might return that it found 1 task, 1 memory, 1 idea, 1 agent assignment, and 2 code chunks in the given conversation <sup>117</sup> – effectively turning an unstructured discussion into structured data.

This conversational interface is crucial for fluid human-AI collaboration. It means users can *talk* to Todozi or an AI integrated with Todozi in a relatively natural way, and the system will update itself accordingly. An AI assistant (like a specialized GPT-based bot) integrated with Todozi can output these tags as it converses. For example, a user might discuss a plan with the AI, and the AI responds with a summary and embedded tasks. The `<todozi>` tag format for tasks expects the fields in a specific order (action; time; priority; project; status; plus optional `assignee=...; tags=...; dependencies=...; context_notes=...; progress=...`), which was designed to ensure completeness <sup>118</sup> <sup>119</sup>. The parsing functions (like `parse_todozi_format` in Rust) validate and create Task objects from these strings <sup>120</sup> <sup>121</sup>. If a field is missing or formatted incorrectly, the parser errors out, ensuring that the conversation produces valid data or the AI learns to output the correct format. This creates a feedback loop: the AI must adhere to a structured output to have its suggestions accepted, which in turn trains it to think in structured terms about tasks.

From a collaboration standpoint, this is powerful. It allows a workflow such as: a human says, "We should also consider security audits." The AI agent can reply, "I will add a task for a security audit in two weeks," and actually include  
`<todozi>Create security audit plan; 2 weeks; medium; security; todo; assignee=agent:tester; tags=compliance</todozi>` in its message. Todozi will then add that as a real task assigned to the tester agent, without the user needing to manually input it. The conversation and the task database stay in sync. Similarly, if the user shares a piece of information like "`<memory>...</memory>`" during a chat (perhaps something they remembered or decided on the fly), it's immediately stored and available to the AI later. This **shared context** ensures the AI is always up-to-date with decisions and information that may arise informally.

Another aspect of collaboration is the **Agent System Integration**. We touched on agent assignments in tasks; Todozi also keeps an "agent registry" (a list of Agent definitions, likely including their name, role, and maybe capabilities) <sup>15</sup> <sup>122</sup>. The CLI allows managing agents (`todozi agent list/show/create/update`) <sup>122</sup>, and the server exposes agent info via endpoints (`GET /agents`) <sup>123</sup>. Agents in Todozi

might have statuses or other properties (the text references agent status tracking and dynamic assignment) <sup>124</sup>, which suggests the system can monitor whether an agent is actively working, idle, or needs input. Human users can see agent assignments in task listings (e.g., tasks prefixed with `agent:` in the assignee field), and they can manually reassign or intervene if needed (`todozi update <id> --assignee human`) to take something over from an AI, for instance). This fosters a cooperative environment where AI is part of the team: visible and manageable, not hidden.

**Use Case – AI Pair Programming:** Imagine a scenario of pair programming between a human developer and an AI agent using Todozi. They communicate in a chat interface. The human says, "Let's break down the login feature." The AI responds with a plan: it outputs a few `<chunk>` entries for the code structure (project, module, method chunks for "User model", "Login API", etc.) and `<todozi>` entries for the high-level tasks (e.g., "Design authentication schema", "Implement OAuth2 login flow; 6 hours; high; backend; todo; assignee=human; tags=auth"). Todozi processes this message: new tasks are created for design and implementation (one assigned to the human, one to maybe a coder agent if the AI deferred some coding), code chunk placeholders are set up for the project. The human can then start coding the user model (one chunk) and mark it done in Todozi, while the AI concurrently starts working on another chunk (perhaps writing a test, as a tester agent). They continuously sync via the Todozi system – if the human finds an issue, they could add a memory "`<memory>User table needs indexing for email field; because of expected search load; importance: medium; long term</memory>`" which the AI will later recall when optimizing the database. The AI might periodically use `recommend_similar` on the tasks the human completed to suggest next tasks ("Since you finished login, maybe set up password reset?"). All of this is facilitated by the structured sharing of information. The human and AI essentially maintain a shared to-do list and knowledge base through Todozi's collaboration features, improving coordination and reducing miscommunication.

**Privacy and Control:** It's worth noting that features like secret memories and share levels give fine control over what data the AI can see or what it can expose. In a team setting, a human user might mark certain memories as AI-only so the AI can use them but not reveal them verbatim to others (useful for things like storing an API key or a private note). Conversely, if the AI generates an idea that is not yet ready for humans, it could keep it private until refined. This encourages an environment where AI can have an "internal monologue" or scratch space (in a controlled way) while still eventually converging all important info into the shared space.

In conclusion, Todozi's human-AI collaboration suite – memories, ideas, conversational parsing, and agent assignments – creates a synergistic workflow. The AI is given tools to **remember context** (so it doesn't ask the same questions or forget preferences), to **propose ideas** (enriching human creativity), and to **take initiative in structuring work** (by parsing and generating tasks in conversation). The human users, on the other hand, gain transparency and control: they see what the AI has noted (memory), what it suggests (ideas/tasks), and can correct or guide it by editing those artifacts using familiar CLI or UI operations. The result is a collaborative loop where the AI's contributions are tangible and editable, just like any team member's, and the **conversation itself becomes a productive interface to manage the project**.

## Terminal User Interface (TUI)

While Todozi can be fully operated via CLI commands and automated through its API, the optional **Terminal User Interface (TUI)** provides a rich, interactive experience that is especially useful for local use by developers. The TUI is a curses-like text interface (built with Rust's `ratatui` library) that visualizes tasks,

analytics, and AI insights in real time within the terminal <sup>125</sup> <sup>126</sup>. This interface is designed to improve user experience by presenting information in panels and interactive elements, which is faster for scanning and decision-making than issuing many separate CLI commands.

**Launching and Architecture:** The TUI is an optional component that can be enabled at compile-time (via a Rust feature flag `tui`). When Todozi is built with this feature, the `todozi tui` command becomes available to launch the interactive UI <sup>127</sup>. Internally, `main.rs` detects the feature and will start `TuiService::run()` instead of the normal CLI loop if the user invokes the TUI mode <sup>127</sup>. The TUI connects to the same underlying `TodoziHandler` and storage as the CLI, meaning it has full access to tasks, projects, agents, etc., and can invoke the same functions but in response to UI events rather than textual commands <sup>128</sup> <sup>129</sup>. There's even a notion of running the TUI in "remote" mode as a client to a Todozi server, which suggests the TUI can connect to a running Todozi service over the network to manage tasks on a remote machine <sup>130</sup> (in such a case, it would use the REST API under the hood to fetch and update data, acting like a graphical client but in the terminal).

**Layout and Panels:** The TUI typically divides the terminal screen into multiple sections (panels) to display different types of information simultaneously. Based on the documentation and code, we can infer some of the panels: - A **Task List Panel** showing active tasks (possibly filtered by project or status). This panel likely lists each task with key fields (ID, description, status, assignee) in a tabular form. Users can navigate up/down to select a task. - A **Task Detail / Editor Panel** that, when a task is selected, shows detailed information about it (full description, all metadata fields, maybe history of updates). This panel might also allow editing: the TUI has an editor state where the user can press a key to edit a field of the task (like change priority or add a context note) and then save changes <sup>131</sup> <sup>132</sup>. In the code snippet we saw, the TUI handles user input for editing each field (mapping typed input to the `editor.current_task` fields such as Priority, Status, Assignee, etc.) <sup>133</sup> <sup>134</sup>. This indicates a form-based editing mode in the UI. - An **Analytics/Visualization Panel** for AI-driven insights. The documentation specifically mentions that the TUI can visualize *similarity graphs, cluster overviews, and drift reports* <sup>135</sup>. Therefore, there might be a panel where the user can toggle to see a graph of related tasks (maybe a network graph drawn with text/ASCII, or a list of similar tasks to the currently selected one), and a panel showing clusters (perhaps listing cluster topics or a dendrogram-like view), and one for drift (like highlighting if the current task's embedding drifted). These could be separate modal views or toggled sections. - A **Charts Panel** illustrating project metrics over time. The code from `tui.rs` shows rendering of two charts: "📈 Task Progress Over Time" and "📊 Activity Over Time" <sup>136</sup> <sup>137</sup>. The *Task Progress Over Time* chart likely plots how task completion percentage evolves over time (maybe each X axis tick is a time unit and Y is % complete or number of tasks completed). The *Activity Over Time* chart appears to plot counts of tasks, ideas, and memories over time (the code was generating sine and cosine waves for demo, but presumably in real usage it would use actual data like how many tasks were created or completed in each time interval) <sup>138</sup> <sup>139</sup>. These charts turn the raw data into visual trends, which is useful for project monitoring (e.g., are we accelerating or slowing down in completing tasks? Did idea generation spike recently?). The TUI draws these with labeled axes "Time" and "Progress %" for the first chart, and similar for the activity chart <sup>140</sup> <sup>141</sup>. - Possibly an **Agent/Queue Panel**: If Todozi has a concept of an agent queue or current agent activities, a panel could list what each agent is doing, or list pending agent actions. There is mention of a `queue` system (with commands like `todozi queue list --status active`) in the CLI context <sup>142</sup>, which might correspond to tasks that are scheduled or being processed by agents. The TUI could present that in a user-friendly way.

The TUI likely uses keyboard shortcuts to navigate between panels and perform actions. Common keys might be j/k or arrow keys to move in lists, Enter to select or drill down, and specific letters for actions (e.g.,

"E" to edit, "A" to add a task, etc.). Although we don't have a direct list of controls in the docs we saw, this is typical of such interfaces. The TUI is essentially harnessing the CLI functions behind the scenes, but giving a continuously updated view (e.g., auto-refreshing the task list when tasks change, showing a live clock or counters).

**Integration with AI Features:** The TUI is not just a static view of tasks; it's built to highlight the AI-augmented aspects of Todozi. For example, when similarity search or cluster features are available, the TUI likely provides an interactive way to use them. Perhaps a user can press `/` to open a search bar and type a query; the TUI would then display search results (across tasks, ideas, memories) in a pane, using the semantic search under the hood. Or selecting a task might show "Similar Tasks" in a side pane with their similarity scores. The mention of drift reports being viewable suggests that if a task changed a lot, the UI could indicate that (maybe an icon or color coding tasks that have high drift). There might also be visual cues such as color highlights for tasks assigned to AI vs human, or priority indicated by color (urgent tasks in red, etc.). The code suggests a color scheme is defined (`self.display_config.color_scheme`) with primary, text, warning colors etc. being used in UI elements <sup>140</sup> <sup>143</sup>.

Another integration point: The TUI can act as a client to the Todozi REST API in remote mode <sup>130</sup>. This means you could run the TUI on your local machine but connect to a Todozi server running elsewhere (for example, a central Todozi server for a team). This is quite useful: it provides a full graphical interface without needing a web browser, good for developers who live in the terminal. In this mode, every action the TUI takes would call an HTTP endpoint (e.g., when you mark a task done in the TUI, it might perform an HTTP PUT to `/tasks/{id}` on the server). The documentation notes that errors from handlers are mapped to HTTP status codes <sup>144</sup>, implying the TUI expects proper responses from the server and can show error messages accordingly.

**Efficiency and User Experience:** The TUI's goal is to make local interaction efficient. For instance, you can scroll through tasks quickly rather than running multiple `todozi list` commands with different filters. Bulk editing might be easier (if the TUI supports multi-select or quick edits). The visual charts and graphs condense a lot of information (progress, activity, relationships) that would otherwise require running separate analysis commands or mentally aggregating. This aligns with the idea that the TUI is for an *AI-augmented project dashboard*. You can imagine leaving the TUI open on one monitor, updating in real time as agents complete tasks and new ones come in, giving the team immediate feedback.

From the development perspective, the TUI is feature-gated and doesn't impact the core logic if not enabled, which keeps the core lightweight for headless or server use <sup>145</sup>. But when enabled, it leverages the same `TodoziHandler` and `Storage` as everything else, ensuring consistency. The UI code interfaces with these through the `TuiService` which likely translates UI events to calls like `handler.add_task(...)`, `handler.list_tasks(...)`, etc. This separation also means the TUI could be improved or replaced without altering the business logic.

**Example Walk-through:** Suppose a user starts Todozi in TUI mode by running `todozi tui` <sup>146</sup>. They are greeted with a screen showing their default project's active tasks in a list. At the top, maybe summary stats: "5 Tasks (3 Todo, 1 InProgress, 1 Done), 2 Memories, 1 Idea" – giving a quick overview. The user uses arrow keys to select a task "Implement chunking system". The right panel displays details: description, created date, etc., and maybe an ASCII chart of its progress (if sub-tasks or chunks existed). The bottom of the screen might show hints or menu: e.g.,

[A] Add Task	[E] Edit Task	[M] Memories	[I] Ideas	[S] Search	[Q] Quit
--------------	---------------	--------------	-----------	------------	----------

The user

presses **M** to open the memory list – the task list panel is replaced by a memory list. They see a memory like "Client prefers iterative dev (long-term, high importance)". Pressing Enter on it shows the full memory details in the detail panel. They press **I** to switch to ideas, etc. At some point, the user wants to see analytics, so maybe they press **T** for toggle charts – the screen splits to show the "Task Progress Over Time" line chart, indicating perhaps that task completion went from 0 to 50% over the last week, and an "Activity Over Time" chart showing that lots of ideas were added yesterday (spike in idea count). The user could also inspect AI-specific info: perhaps pressing **R** for related tasks pops up a window listing tasks similar to the current one, with similarity scores. If they press **C** for clusters, the UI could show a list of cluster names or a cluster tree. They might see cluster #3 labeled "Auth / Login (5 items)" – selecting it shows those 5 tasks/ideas that are semantically grouped.

This kind of TUI is invaluable for quickly traversing the knowledge base that Todozi manages, and it's particularly appealing to developers who prefer terminal environments. It condenses the collaborative, AI-infused data (tasks, memory, ideas) into a navigable interface that feels like a live dashboard. Additionally, because it can integrate with remote mode, a developer can SSH into a server running Todozi and use the TUI from there, managing tasks on a headless server with ease.

**Real-World Use Case for TUI:** A startup team might run the Todozi server on a cloud VM so that all AI agents and storage run centrally. Each developer runs `todozi tui --server <address>` from their own machine to connect. They get the same view of the project and can see in real-time as tasks move (for example, an AI agent marks a task done — it might blink or change color in their TUI, or a new idea appears when someone else added it). They can still use their keyboard to quickly add or modify tasks, which sends updates to the server. This merges the benefits of a web-based project management tool (multi-user, live updates) with the speed and scriptability of terminal applications. And since Todozi's TUI can render advanced visualizations (progress charts, etc.), team members can monitor AI agent activity (like a chart of tasks completed by AI vs humans over time) easily. It's like having a power-user interface for an AI-driven JIRA/Trello, all in a terminal window.

---

## Conclusion and Future Perspectives

Todozi represents a new breed of task management system that is **"AI-first" in design, yet human-friendly in practice**. Through this detailed exploration, we have seen how Todozi's architecture weaves together structured project management, semantic AI capabilities, continuous model improvement, and collaboration features into a cohesive tool for productivity.

To recap the key points: - **AI-Integrated Task Management:** Todozi's robust task model and agent assignment framework enable tasks to be naturally shared between humans and AI agents, with clear metadata (priority, status, dependencies, etc.) that both can interpret <sup>3</sup>. Hierarchical decomposition (like the chunking system) allows tackling large problems in manageable parts, aligning with how AI processes context in chunks <sup>26</sup> <sup>147</sup>. - **Embedding & Semantic Search:** By converting all content into semantic vectors, Todozi unlocks intelligent features like meaning-based search, automatic clustering of related work, and AI recommendations for next tasks <sup>5</sup> <sup>53</sup>. The embedding service is high-performance and extensible, supporting model swaps, multi-model comparisons, and ongoing quality checks <sup>71</sup> <sup>56</sup>. - **Model Training Pipeline:** The system doesn't rely on static AI models; it provides mechanisms to learn from usage. Data logging and export allow fine-tuning models to specific project domains, while versioning and



validation ensure that these AI improvements integrate smoothly without loss of integrity <sup>81</sup> <sup>7</sup>. This means Todozi's AI can get smarter and more customized over time – a critical factor for real-world deployments where one size doesn't fit all. - **Human-AI Collaboration:** Todozi treats the AI as an equal stakeholder in the project space. Features like the memory system and idea logs capture context and creativity that inform both human decisions and AI reasoning <sup>9</sup> <sup>148</sup>. The conversational interface transforms how users interact with management tools – instead of filling forms, they can simply discuss plans and let Todozi parse out tasks and notes <sup>100</sup> <sup>10</sup>. This lowers the friction of maintaining a project plan and ensures important details from discussions are not lost. - **Terminal UI:** For power-users and developers, Todozi's TUI offers a fast, immersive view into all this data. It blends traditional task lists with AI-driven analytics (like progress charts and similarity graphs) <sup>140</sup> <sup>135</sup>, providing an at-a-glance understanding of project status and AI contributions. The TUI exemplifies Todozi's ethos of combining the new (AI insights) with the familiar (straightforward UI/UX) to maximize adoption.

**Extensibility and Integration:** Todozi's modular architecture (Rust library + CLI + optional server) makes it highly extensible. New features can be added as new command groups or modules (for example, adding a `calendar` integration or a `time tracking` module) without disrupting existing functionality <sup>149</sup>. The presence of a REST API with full coverage of core operations (tasks, projects, memories, ideas, agents, chat processing) means Todozi can be integrated into other systems or UIs easily <sup>150</sup> <sup>151</sup>. For instance, a web dashboard or a VS Code plugin could use the API to present Todozi data in those contexts. The API also enables multi-language client libraries – the documentation even provided sample JavaScript and Python client usage <sup>152</sup> <sup>153</sup>, showing how any environment can drive Todozi (e.g., a Python script could automatically add tasks from a bug tracker, or a Node.js app could feed in chat from a Slack channel to Todozi for processing).

**Real-World Impact:** The capabilities described are not just theoretical – they address common pain points in knowledge work and software projects. Teams often struggle to keep documentation, to-do lists, and decision logs up to date because it's tedious. Todozi, with AI assistance, automates much of that: it can parse meeting notes (via chat processing) into updates, suggest tasks that might be missing, and ensure nothing falls through the cracks by semantic linking. It effectively serves as an AI project manager alongside the human project manager. Early adopters could be software dev teams (using it to integrate with code generation and testing agents), research groups (managing literature reviews and experiments, with the memory system capturing key findings), or even individuals managing complex projects with lots of moving parts (who would benefit from the AI suggesting next steps and remembering details).

**AI-First Design Benefits:** By structuring data for AI, Todozi also sets the stage for advanced AI behaviors. For example, one could plug in a large language model agent that reads the entire task list and memory and generates a project summary or risk report. Since everything is structured (and can be serialized to JSON or HLX), feeding it into an AI and getting outputs that Todozi can parse is straightforward. This two-way synergy (Todozi <-> AI) is a glimpse of how future productivity tools will operate, and Todozi is a concrete realization of that vision.

**Future Directions:** The documentation hints at possible future enhancements that were beyond the initial implementation <sup>154</sup>. These include a web interface (for broader user adoption), plugin system, and deeper analytics. A particularly interesting future feature is **streaming embeddings** (deferred in the current version) <sup>155</sup> – this could allow real-time updating of embeddings as text streams in (imagine as you type a task, the system could live-update similar tasks). Another is **graphical visualization dashboards** for the embedding space <sup>156</sup>, which could bring even more insights (e.g., a 2D projection of all tasks by similarity).

With the architecture in place, such additions would enrich the platform but are not required for core functionality thanks to the thoughtfully designed foundation.

In conclusion, Todozi exemplifies an **AI-augmented task management system** that doesn't just manage what you tell it, but actively helps you manage better. It merges the reliability and organization of traditional project management with the intelligence and adaptability of modern AI. For AI researchers, Todozi provides a playground for multi-agent collaboration and continual learning in a real-use case domain. For developers, it offers a powerful tool to offload routine planning drudgery to AI while retaining oversight and control. The extensible, open design invites experimentation – whether integrating new AI models, adding custom agents, or connecting Todozi to external data sources. As AI becomes an integral part of daily workflows, systems like Todozi will likely become the norm, serving as the connective tissue between human creativity and machine efficiency. This whitepaper has detailed how Todozi's system achieves this integration today, setting a high bar for what it means to be “AI-first” in the realm of task and project management.

Sources:

- Todozi Implementation & Documentation Files 1 3 5 10 etc. (Provided in user documentation)

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30

31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90

91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 122

123 124 125 126 127 128 129 130 135 142 144 145 146 147 148 149 150 151 152 153 154 155 156 all-tdz.md

file:///file\_00000000740061f5b08286d9d20607ed

118 119 120 121 todozi.rs

file:///file\_0000000066c061f7b844a16a7e35068d

131 132 133 134 136 137 138 139 140 141 143 tui.rs

file:///file\_00000000bfd061f7af9cd5ceb0a4e556