

# Local AI Architecture Reference

Architectural reference for the Local AI capability in FrameworkX 10.1.5: the two-path execution model, code organization, master-gate semantics, transcript cache mechanics, tool dispatch loop, and integration points with the rest of the platform.

[Technical Reference](#) [Platform Architecture Reference](#) [Local AI Architecture Reference](#)

Version 10.1.5+

## Audience and scope

This page is for system integrators, custom-driver authors, OEM partners, and engineers who need to reason about how Local AI is composed inside the FrameworkX runtime: which assembly hosts what, which thread runs which work, where gates apply, and which contracts must be honoured when extending or wrapping Local AI.

End users configuring Local AI from the Designer or wiring a chat panel from a Display do not need this page. See [Local AI](#) and the user-facing reference children there.

## Capability summary

Local AI exposes one platform capability — calling an OpenAI-compatible LLM from inside a FrameworkX solution — through two architecturally independent consumer paths:

Path	Surface	Caller context	Stateful?	Tools?	Hooks?
<b>Cached</b>	ChatRequest Display action	Operator click ServiceClient TCP Kernel TCPServer dispatch Local AI	Per-Display-panel transcript cache	Yes (UNS / Alarm / Historian / custom)	Yes (OnBeforeChat / On AfterChatReply)
<b>Atomic</b>	TK.AIExecute script API	Server-domain script (Server.Class, Script.Task, alarm callback, report generator) Toolkit static	Stateless — every call independent	No (single POST, no tools)	No

Both paths POST to the same OpenAI-compatible endpoint, configured once per solution on three columns of `SolutionSettings`. Both paths return the same reply envelope (see [Local AI Reply Envelope Schema](#)).

## Two-path model — what is shared and what is divergent

The two paths are deliberately independent code paths sharing a small set of explicit anchors. Code duplication is welcomed where the alternative would be a shared helper that couples them.

### Shared by architecture (legitimate coupling)

- **The endpoint configuration** — both paths read `SolutionSettings.ModelSettings` JSON via the same defensive accessors on the Local AI service singleton.
- **The master kill-switch** — both paths check `SolutionSettings.ModelEnabled` at entry and short-circuit identically when off.
- **The SecuritySecrets resolver** — both paths invoke `Security.CheckAndResolveSecretCommConfiguration` on resolved auth strings; this is a Kernel-side platform service, not a path-specific helper.
- **The reply envelope schema** — both paths produce the same five-field JSON shape (`text, status, toolTrace, latencyMs, warnings`).
- **The auth multi-line decoder** — both paths call the platform's pre-existing `ReportWebData.DecodeAuthorizationHeader` static, which is the single source of truth for the `None / BearerToken / BasicAuth / CustomAuth` wire format.

### Divergent by design

Concern	Cached path	Atomic path
Implementation class	Sealed orchestrator owning the per-Display-panel transcript	Separate sealed class — no shared methods, no shared state
HTTP-POST helper	Inline in the cached orchestrator	Inline in the atomic class — separate copy, even if today the bodies look similar, so divergence as features grow stays isolated
Master tool-surface bit ( <code>ModelOptions.0x02</code> )	Required — off <code>status="disabled"</code>	Not consulted — atomic has no tools to gate
Per-category tool bits ( <code>0x04 - 0x20</code> )	AND-ed against master to assemble the tool catalog	Not consulted
Hook chain	Raises <code>OnBeforeChat</code> and <code>OnAfterChatReply</code>	Does NOT raise hooks — chat-attached behaviour does not apply
Transcript cache	<code>PerClientInfo.Guid</code> , lazy login-reset, <code>0x80 EnableChatHistory</code> sub-gate	None

Tool dispatch loop	Bounded by per-turn cap and 60-second wall-clock budget	Single POST, no loop
Caller wire signature	Receives (ObjectServer, clientGuid, userName, queryJson)	Receives (ObjectServer, queryJson)

## Code organization

### Assembly home and the namespace-contribution rule

Local AI is implemented as a sub-namespace inside `T.Toolkit.dll` — the `T.Toolkit.LocalAI` namespace, with sources under `FS/Toolkit/LocalAI/`. It is NOT a peer assembly to `T.Toolkit.dll` and it is NOT inside `T.Modules.dll`.

The placement is governed by a platform rule: *code lives in T.Modules.dll if and only if it contributes to the runtime Object Model* — Tag namespaces, `eObjType` registration, `ListObj/RunObj` entries, user-visible runtime addressability. Local AI has zero Object Model contribution: `TK.AIExecute` (a Toolkit static) and the `Display ChatRequest` action both *consume* the Object Model — they read tags, dispatch `runtime_*` tools — but neither *contributes* a Tag namespace or a runtime root. Therefore Local AI lives as a sub-namespace inside `T.Toolkit.dll` rather than as its own peer assembly or as a runtime Module.

This explains why there is no `AI` namespace in the runtime Object Model in 10.1.5 — Local AI is reachable only as `TK.AIExecute(...)` from script and as the `ChatRequest Display` action. A future `@Solution.AI.*` binding surface is post-10.1.5 work.

### Build position and dependencies

Assembly	Local AI sources	Depends on	Consumed by
<code>T.Toolkit.dll</code> ( <code>T.Toolkit.LocalAI</code> sub-namespace)	<code>FS/Toolkit/LocalAI/*.cs</code> , included in <code>T.Toolkit.csproj</code> via the implicit SDK glob — no explicit <code>&lt;Compile Include&gt;</code> entries	<code>T.Kernel.dll</code> , <code>T.Library.dll</code> , <code>T.ProjectServer.dll</code> (already required by <code>T.Toolkit.dll</code> )	<code>TK_AI.cs</code> (atomic-path entry, intra-assembly direct call), Kernel-side TCP service host (cached-path entry, cross-assembly), <code>T.Designer.UI</code> (defaults class)

The Local AI sources are included in `T.Toolkit.csproj` via the implicit SDK glob under `FS/Toolkit/LocalAI/` — no separate project, no project reference. The Toolkit consumer (`TK_AI.cs`) calls the Local AI service via an intra-assembly direct call — same DLL, same compile unit, no reflection-by-name on this binding.

### Bind sites

Two callers reach the Local AI service entry methods. Both are compile-time-checked under the prescribed code organization, so signature changes surface as build errors rather than runtime nulls.

Caller	Binding	Calls
Kernel-side TCP service handler (cached path)	Cross-assembly compile-time reference from the Kernel-side service host to <code>T.Toolkit.dll</code> (the cached-path service lives in the <code>T.Toolkit.LocalAI</code> sub-namespace)	<code>ExecuteChatAsync</code> ( <code>ObjectServer</code> , <code>clientGuid</code> , <code>userName</code> , <code>queryJson</code> )
Toolkit static <code>TK.AIExecute</code> (atomic path)	Intra-assembly direct call inside <code>T.Toolkit.dll</code> ( <code>TK_AI.cs</code> and the <code>T.Toolkit.LocalAI</code> service compile into the same DLL)	<code>ExecuteAtomicAsync</code> ( <code>ObjectServer</code> , <code>queryJson</code> )

### Type visibility

The Local AI service singleton, the cached-path orchestrator, and the atomic-path class are all `internal sealed`. The single `public static` surface is the defaults class, which exposes the recommended endpoint constants (URL, model name, etc.) for the Designer-side configuration dialog to seed its 5-field editor.

## Process and thread model

### Where Local AI runs

Local AI HTTP traffic is server-side only. Both paths execute their LLM POST inside the `TServer` process; clients never directly contact the LLM endpoint.

Caller surface	Process where LLM POST happens	Thread
<code>ChatRequest Display</code> action	<code>TServer</code> (TCP service handler dispatch)	<code>TCPServer</code> dispatch thread for the inbound <code>ServiceClient</code> call

TK.AIExecute from Server. Class method	TServer	Caller's thread (script execution thread); sync wrapper unwraps async via AsyncHelpers.RunSync
TK.AIExecute from Server. Task	ScriptTaskServer.exe (a server-domain process)	Task execution thread
TK.AIExecute from Report generator	ReportServer.exe	Report-rendering thread

## Why no runtime "view-client" gate

The cached path is reachable only through the Kernel-side TCP service handler — by construction, that handler runs in the TServer process. A misrouted client-side caller cannot land in `ExecuteChatAsync`.

The atomic path is reachable from a Toolkit static, so *technically* a Display CodeBehind script (which runs at the client) could call `TK.AIExecute`. The platform does not block this with a runtime gate, because three platform-level safety nets already defend against degraded execution at the client:

- Secret resolution short-circuits off-server.** `Security.CheckAndResolveSecretCommConfiguration` returns the unresolved literal / `secret:<Name>` token when not running server-side. The LLM endpoint then receives the literal token as the credential and rejects with 401.
- solutionSettings reads route through ObjServer.DB.** Some client-side configurations have a stripped-config DB; reads fall back to defaults targeting `localhost:11434` — unreachable from a remote browser.
- Master kill-switch precedes everything.** When `ModelEnabled` is off, the call short-circuits with `status="disabled"` before any HTTP work.

Together these mean a misrouted call from a thin client never silently succeeds with leaked credentials — it returns a normal HTTP-error envelope or the master-gate disabled envelope. No runtime view-client predicate is required.

## Master gates — application order

Both paths apply gates in a fixed sequence at entry to the service method. The order is non-negotiable.

Order	Gate	Cached path	Atomic path	Outcome when gate is closed
1	<code>SolutionSettings.ModelEnabled = true</code>	Yes	Yes	<code>status="disabled", latencyMs=0, no HTTP traffic</code>
2	<code>SolutionSettings.ModelOptions</code> bit 0x02 (EnableRuntimeMCP)	Yes	Skipped	Cached: <code>status="disabled"</code> ; atomic: not applicable
3	Per-category sub-bits (0x04–0x20)	AND-ed with master to build tool catalog	Skipped	Tool category absent from the catalog the LLM sees

The asymmetry on gate 2 is intentional. `ModelOptions` exists to gate the LLM's tool surface; the atomic path exposes no tools, so the bit is meaningless to it. Gating the atomic path on 0x02 would prevent customers from using `TK.AIExecute` for pure-language tasks (translation, summary, rephrase) while reserving tool-equipped chat for selected operator panels — a legitimate posture.

## Cached-path mechanics

### Per-Display-panel transcript cache

The cached orchestrator maintains a per-solution `ConcurrentDictionary<string, ChatSession>` keyed by `ClientInfo.Guid` — the platform's per-TCP-connection identifier (the same key shape `ModuleHistorian` and `ServiceSyncObjects` use for their own per-connection bookkeeping).

Property	Behaviour
Cache key	<code>ClientInfo.Guid</code> from the inbound TCP service call's thread-local context. NOT exposed to the runtime <code>Client.*</code> namespace.
Cache participation	Gated by <code>SolutionSettings.ModelOptions</code> bit 0x80 (EnableChatHistory). Default ON in new 10.1.5 solutions. When OFF, the orchestrator builds a local <code>SessionState</code> per call without entering the cache.
Login reset	Lazy. The orchestrator records the last-seen <code>UserName</code> per cache entry; when the next chat turn arrives with a different <code>UserName</code> on the same <code>ClientInfo.Guid</code> , the transcript clears before processing the new turn.
Transcript truncation	Per-cache-entry, truncated to a fixed maximum number of messages calibrated for ~8B-class context windows. Older messages drop oldest-first.
Cache lifetime	Per-process. Restarting TServer clears all transcripts. There is no persistence to disk in 10.1.5.

## Tool dispatch loop

When the master tool bit is on and at least one category sub-bit is on, each chat turn assembles a tool catalog and may enter a tool-dispatch loop:

1. POST `messages + tools` to the LLM endpoint.
2. Inspect the assistant message. If it contains `tool_calls`, dispatch each tool in-process against the live `ObjectServer` (no HTTP loopback to `RuntimeMCP` — sub-millisecond direct dispatch).
3. Append tool-result messages to the running message list. Loop back to step 1.
4. Terminal exit when the assistant message is plain content (no `tool_calls`).

The loop is bounded by two limits:

- **Per-turn dispatch cap** — total tool calls per turn. `MaxToolDispatchesPerTurn = 5` in **10.1.5**. Conservative cap chosen for the wall-clock budget below: at typical CPU-LLM latencies of 3–15 s per round-trip, 5 dispatches plus the terminal POST (worst case ~6 LLM POSTs) keeps the turn inside 60 seconds. On GPU-class hardware (~0.5–2.5 s per POST) the cap is not the binding constraint — the wall-clock budget is. When the cap trips, the next iteration POSTs WITHOUT new tool results so the LLM gets one final chance to compose a terminal reply, and returns `status="truncated"` if it does not.
- **Wall-clock budget** — 60 seconds covering the entire loop (LLM POSTs + in-process tool dispatches combined).

## Hook chain

`OnBeforeChat` and `OnAfterChatReply` are event `Func<string, Task<string>>` on the per-`ObjectServer` Local AI service singleton. Multicast; throwing handlers caught and logged into the reply's `warnings[]` with handler identity; chain proceeds with the un-mutated input. Hook failures do NOT flip the envelope's `status` away from `ok` — that's reserved for chain-aborting failures (LLM HTTP error, malformed query, etc.).

See [Local AI Developer Reference](#) for the reflection-attach pattern and worked hook examples.

## Atomic-path mechanics

The atomic path is deliberately simple. After the master `ModelEnabled` gate:

1. Parse query JSON (plain text wrap or structured form with `system / user / context / metadata` fields).
2. Read fresh `ModelSettings` via the same defensive accessors as the cached path.
3. Resolve `/secret:<Name>` tokens via the platform's secret resolver.
4. Resolve `{Tag.X}` expressions in URL via the platform's expression resolver.
5. Decode auth multi-line via `ReportWebData.DecodeAuthorizationHeader`.
6. Build a single OpenAI-compatible `messages[]` array — no transcript prefix, no prior turns.
7. POST with empty `tools[]`, 60-second wall-clock `CancellationTokenSource`, no loop.
8. Read `choices[0].message.content`. Build the SPEC envelope (`toolTrace=[]` always for atomic).

No hooks. No transcript. No tool dispatch. No retry. Every step lives in the atomic class; nothing calls into the cached orchestrator.

## Configuration topology

### Three columns on `SolutionSettings`

Local AI uses three pre-existing scaffolding columns. No schema migration; no per-solution row creation; no AI-specific table.

Column	Type	Role
<code>ModelEnabled</code>	Boolean	Master kill-switch. Off short-circuits both paths.
<code>ModelSettings</code>	String (JSON blob)	Five-key endpoint configuration: URL, Name, Authorization, Headers, Info.
<code>ModelOptions</code>	Int (bitmask)	Master tool bit and per-category sub-bits. Shared bitmask with the AI Runtime Connector and the AI Designer connector.

### Defensive accessors

Every Local AI call re-parses `ModelSettings` JSON on every read. There is no caching layer with invalidation. The parse cost is negligible compared to the LLM round-trip; the alternative (cached parsed values + invalidation on column change) introduces a stale-cache bug surface that buys nothing measurable. Reads of NULL / empty / malformed JSON / missing keys / empty values transparently fall back to the defaults class.

### Forward-compatibility on unknown keys

The Designer's 5-field editor preserves unknown JSON keys via a capture/replay mechanism, so future SPEC extensions to `ModelSettings` (e.g., a `Provider` dialect enum) survive an open/OK round-trip in the editor without being stripped.

## SecuritySecrets integration

Local AI uses the same SecuritySecrets pattern that protocol drivers, the WebData connector, the UNS tag-provider service, and the Devices module use for credentials. There is no AI-specific secret mechanism.

Aspect	How Local AI uses it
Reference syntax	<code>/secret:&lt;Name&gt;</code> tokens embedded in URL, Authorization, or Headers fields of ModelSettings
Resolver	<code>ObjectServer.Security.CheckAndResolveSecretCommConfiguration(string)</code> — Kernel-side platform service
Off-server behaviour	Resolver short-circuits when not <code>RunningAtServer</code> ; literal token survives. Calling code falls through to a normal HTTP-error envelope.
Encryption	<code>SecretValue</code> is encrypted at rest. Field is write-only via the Designer Security UI; redacted on read-back.

For the customer-facing walk-through and worked examples, see [SecuritySecrets Authentication for Local AI](#).

## Reply envelope contract

The reply envelope is a strict five-field shape, identical for both paths and identical across success / error / disabled / truncated outcomes:

Error rendering macro 'code': Invalid value specified for parameter 'com.atlassian.confluence.ext.code.render.InvalidValueException'

```
{
  "text":      "<answer text or '' on non-ok status>",
  "status":    "ok | error | disabled | truncated",
  "toolTrace": [ ... ],
  "latencyMs": 1247,
  "warnings":  [ ... ]
}
```

Architectural guarantees:

- **Always parseable.** Every envelope is well-formed JSON. `JsonObject.Parse(reply)` never throws on a Local AI reply.
- **Field stability.** All five fields present on every envelope, every code path. Adding new top-level fields is non-breaking.
- **latencyMs always present** — including on errors. Pre-condition errors (gates) report 0; post-condition errors report actual elapsed time.
- **Never throws.** Network failures, parser failures, gate-off conditions, internal exceptions — all become envelopes with `status="error"` or `status="disabled"` and a populated `warnings[]`.

Full schema and worked examples: [Local AI Reply Envelope Schema](#).

## Integration points with other platform services

Service	How Local AI uses it
ObjectServer (per-host runtime)	Local AI service is per-ObjectServer singleton ( <code>LocalAIService.Get(objSrv)</code> ); accessors route configuration reads through <code>objSrv.DB</code> .
Security (SecuritySecrets resolver)	Token resolution at call time, server-side only.
SolutionSettings (configuration)	Three columns carry all configuration; no AI-specific table.
Kernel TCPService (cached path entry)	Hosts the <code>ChatProcessService</code> handler that the Display ChatRequest action calls via <code>ServiceClient TCP</code> .
Toolkit (atomic path entry)	<code>TK.AIExecute (TK_AI.cs)</code> reaches the Local AI service via an intra-assembly direct call — both <code>TK_AI.cs</code> and the <code>T.Toolkit.dll</code> .
WebData connector	Reuses <code>ReportWebData.DecodeAuthorizationHeader</code> for auth multi-line decoding — single source of truth for the wire format.
RuntimeMCP tool catalog	Cached path's tool catalog mirrors the same <code>runtime_*</code> tool surface that <code>RuntimeMCP.exe</code> exposes to external AI clients — single tool definition, two transports (in-process for cached path, HTTP for external clients).
AsyncHelpers	Both paths use <code>T.Library.tRPC.AsyncHelpers.RunSync</code> for sync-over-async wrapping (NOT raw <code>.Result / .GetAwaiter().GetResult()</code> ) — handles <code>SynchronizationContext</code> correctly and unwraps <code>AggregateException</code> .

## What Local AI deliberately does NOT add

Not added in 10.1.5	Why
---------------------	-----

An AI runtime root in the Object Model	The 12-root taxonomy (Tag, Server, Client, Alarm, Device, etc.) does not gain an AI root. Surface remains <code>TK.AIExecute + ChatRequest</code> . A peer of <code>InfoSolution.Settings</code> exposing AI fields is post-10.1.5.
A <code>T.Modules.AI.TModule</code>	Local AI does not contribute Object Model namespace surface; it lives as a sub-namespace inside <code>T.Toolkit.dll</code> , not as a runtime Module.
New <code>SolutionSettings</code> columns	The three existing columns ( <code>ModelEnabled</code> , <code>ModelSettings</code> , <code>ModelOptions</code> ) carry everything. Schema migration would impose cross-version routing risk for no benefit.
A provider-dialect enum (OpenAI / Ollama / Anthropic)	The dialect is implicit (OpenAI-compatible JSON). Adding a discriminator waits for a real customer needing a non-conforming endpoint.
Streaming output	Sync request/response only. Each call returns one complete envelope when the model finishes.
A bundled LLM runtime or model weights	Customers install Ollama themselves and pull the model from the official registry. Keeps the FX installer lean and avoids redistribution licensing friction.
HTTP loopback to <code>RuntimeMCP</code> for tool dispatch	Tool dispatch is in-process direct against the live <code>ObjectServer</code> — sub-millisecond vs 10–50 ms over-the-wire to localhost. Also guarantees the cached path works in standalone deployments without <code>RuntimeMCPHttp</code> running.
Atomic-path hooks or atomic-path tools	Atomic is a deliberately stateless one-shot. A customer wanting hooks or tools should use the <code>ChatRequest</code> Display action surface.

## Versioning and forward-compatibility expectations

- **Reply envelope schema is LOCKED.** Adding new top-level fields in future versions is non-breaking; existing five fields stay present and stay typed as documented.
- **Configuration topology is LOCKED.** The three-column shape will not change. Future configuration extensions land as new keys inside the `ModelSettings` JSON or as new bits in `ModelOptions`.
- **Two-path independence is LOCKED.** Future features ship in one path or the other (or both, with separate implementations). The two paths will not be merged behind a shared helper.
- **Code organization may evolve.** The `T.Toolkit.LocalAI` sub-namespace home is stable for 10.1.5; later releases may split or rename internal types. The compile-time entry surface (`ExecuteChatAsync`, `ExecuteAtomicAsync`) is signature-stable.

## See also

- [Local AI](#) — the customer-facing capability page (overview + quick start).
- [Local AI Configuration](#) — endpoint, master enable, and bitmask configuration.
- [ChatRequest Action Reference](#) — the operator-chat surface.
- [TK.AIExecute API Reference](#) — the atomic script surface.
- [Local AI Reply Envelope Schema](#) — the uniform reply envelope.
- [SecuritySecrets Authentication for Local AI](#) — credential management.
- [Local AI Developer Reference](#) — hook attachment, custom tools, advanced patterns.

---

## In this section...