## Our client, a global insurance aggregator, scales by adding new

Business challenge

countries. Each partner comes with different APIs, schemas, languages (including non-Latin scripts and right-to-left layouts), and regulatory constraints. Historically, onboarding a single partner took **3-6 months** of cross-functional effort: clarifying requirements, interpreting

partners (carriers, MGAs, regional brokers) across dozens of

sparse and heterogeneous documents, writing adapter code, preparing test data, iterating through compliance checks. Multiplied by hundreds of partners, the cost ballooned and timelines stretched.

Compress time-to-integration from months to weeks

The client needed to:

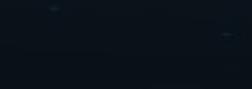
without compromising compliance or auditability; Accommodate variability in partner maturity: from bringing different API protocols (REST, SOAP, etc.) to

a common denominator to handling multiple document types (PDFs, spreadsheets, or even email samples);

Let non-technical partner representatives self-serve 03 in their native language, with a guided, transparent flow; Stay model-agnostic and cost-efficient as the LLM 04

Instinctools' dedicated AI team took on creating a faster and smoother partner onboarding workflow.

landscape evolves.



01

Solution

combines a structured Al adoption process, an orchestration pipeline that thinks before it codes, and model governance from our AI Center of Excellence.

automated tests, finishing with a GitHub pull request. The solution

We delivered a **production-ready**, **UI-first**, **multi-agent system** that turns

partner inputs (documents, answers, samples) into working adapters and

To de-risk the initiative and align on outcomes, we started with the Al adoption

Running an Al

adoption workshop

Conducted business discovery, mapping the core Defined clear success criteria and guardrails with a flows (Quote → Bind → Endorsement → Cancellation → target of ~2 weeks per partner, compile-clean builds, Claims/FNOL), the country-specific nuances, and the minimum test-coverage thresholds, and audit-ready

workshop. Within it, our team:

- compliance checkpoints; Identified process bottlenecks and time sinks, including document triage, field mapping, eligibility rules, error semantics, and under-specified specifications;
- Selecting the right
- artifacts; Produced solution blueprint around a three-phase agentic pipeline (Analyze → Plan → Generate) with human-in-the-loop validation.

02

Before designing the agentic pipeline, our Al engineers tested multiple LLMs using reverse-engineered samples

model

(Opus and Sonnet), Claude Opus 4.1 delivered the most stable and production-ready output, especially when guided by structured prompts and incremental

from existing API adapters to see which fits the client's

needs best. Among GPT, Gemini, Grok, and Anthropic

repeatable framework that evaluates models by code-gen accuracy, context window, speed,

Model governance by our Al Center

Choosing and re-choosing models is integral to

value. Instinctools' very own AI CoE runs a

of Excellence

modality coverage, hosting options, API availability, cost per token, language coverage, and deployment constraints. Our experts continuously benchmark new releases and propose controlled switches (with cost deltas and risk notes), so the client benefits as the market shifts.

Each adapter runs through a 13-step playbook where every

output becomes the next input, while the agent distills

patterns from the prior implementations into a

consolidated guide for what to build. A short intake

Agents don't jump to code. They analyze inputs, plan the work with acceptance checks, then generate code and

pipeline

Building an agentic

tests, looping until the build is clean.

questionnaire captures the business rules the docs miss, and we seed the workspace with curated reference repos so the agent can "look up" proven approaches. With an agentic approach, a large endpoint is typically completed in 2-3 hours for roughly \$50-100 of model spend.

**Ensuring quality** and auditability

(schema), and integration tests, a compile-fix loop, and smoke/HTTP probes that prove the service actually runs. Humans step in only at high-leverage moments, dedicating, on average, up to 20 minutes to polishing the

output.

For Al-powered onboarding to be safe and sound, we had

to ensure compliance with a plethora of security policies

Every step is pinned to hard checks, such as unit, contract

controls into prompts ("out of scope" rules), keeping models focused on what's required and nothing more. And when something slips, those traces collapse time-toroot-cause from hours to minutes, so the next run captures the fix by design.

Under the hood, we trace prompts, tool calls, inputs/

outputs, latencies, and token spend with Langfuse. That

observability makes audits boring in the best way: each

decision and artifact is tied to a PR with change logs and

Following our Alin SDLC practices, we've also baked drift

test evidence, ready for reviewers and regulators.

To check all the boxes, our team zeroed in on:

Regional hosting options to store EU data within the

PII protection measures like field-level masking to

ensure no sensitive data is fed to the model

06

05

by design: Core security and privacy baselines, such as NIST and **OWASP** 

Putting a premium

on AI security controls

500 and the NAIC Insurance Data Security Model Law for the US, the DORA for the EU, the FCA/PRA Operational Resilience for the UK, and the APRA CPS 234 for Australia, among others

regulations by region, such as the NYDFS 23 NYCRR

Insurance and financial-sector cybersecurity

- Data privacy laws, such as CCPA, PIPEDA, GDPR, PDPA, APPI, etc. Al governance regulations, including NIST AI RMF and
- Prompt "guardrails" to prevent the model from going off-scope or fabricating logic Human-in-the-loop reviews of Al artifacts, with all of them linked to a GitHub PR for audit trails

EU, US data within the US, etc.

Role-based data access

**EU AI Act** 

Designing

self-service UI

simple web UI.

Once the agentic approach proved consistent across 10

API adapters, our AI team wrapped the workflow in a

Human expertise still guides quality and regulatory soundness. The agent automates the mechanical steps, delivering a compile-clean adapter with

Progress and transparency. Step-by-step status, logs,

downloadable reports, and validation results.

Partner-facing portal. Secure onboarding wizard and

Postman/XLSX or describe their process in free text.

partner's language and normalizes vocabulary into the

Multilingual by default. The agent converses in the

chat. Partners **drag-and-drop** PDFs/Swagger/

unit tests and CI checks, which used to consume months of execution time.

Human touch

pages if flagged

Provide access keys if required

Launch a quick sanity check for

hallucinations on big reports

Upload docs, remove obviously irrelevant

platform's canonical model.

### **Documentation ingestion** Parses PDFs/Postman collections, extracts relevant sections per step/endpoint

Agent-powered partner onboarding

with a human in the loop

Stage

Partner intake and registration

Contract and reference analysis

Implementation plan Scans repo, identifies existing services/DTOs, Approve/adjust plan if edge cases proposes what to add/change

mines patterns from prior adapters

Builds reports on aggregator API contract;

At the end of each step of the workflow, a developer reviews the output before the model proceeds. Such application of human

Self-service sign-up, chat onboarding in any

judgment where it matters most makes the agentic onboarding approach predictable and trustworthy, not a "black box".

Automated by agents

language, guided file upload

Code generation (per endpoint)	Generates services/DTOs/mappers/tests aligned to contract	
Build and self-fix loop	Compiles project, iterates on compile errors until green, starts app	
Smoke checks and PR	Optional HTTP smoke tests, opens a GitHub PR with artifacts	Review PR, merge
Orchestration and audit	One-click run from UI, step-by-step logs/ metrics (LangFuse)	Monitor runs, rerun if a step stalls
Compliance and escalation	Routes ambiguity/questions via chat, escalates to legal/SME when needed	Answer a short checklist (e.g., coverage limits, exclusions, etc.)
UI SOMEN KHULLI III JUNE III JUNE		
Before	Afte	er
Fragmented documentation in multiple languages		eeks of onboarding time per partner, end-to-end, uding review and PR verned model policy that keeps quality up and token
Slow feedback loops and unclear ownership of quality gates		ts down
Client's software engineers do all the groundwork		ent's engineers take on the human-in-the-loop role while entic AI carries on the integration process step-by-step
<ul> <li>Each adapter is treated as a one-off build</li> </ul>		owledge now compounds: the more adapters are built, easier it is to build the next ones

- **Business value**
- Up to 12× faster partner onboarding 10× decrease in operational costs
- 80–90% less repetitive development work Near-zero rework on recurring issues Full auditability of Al actions for compliance

With the hard parts of partner onboarding automated and a

expand faster into new countries and niches. As the library of

model strategy that evolves with the market, the platform can

reference adapters grows, each subsequent integration benefits

# Multiplier effect

from pattern reuse, shrinking effort even further.

Contact us — and we will estimate your projects costs for free!

Do you have a similar project idea?

**CONTACT US**