

Al-Powered Optimization for Container & White Glove Logistics

XGBoost, GPT, and OCR models working together to streamline high-complexity logistics.





Siliconmint partnered with a U.S. logistics provider to transform their transportation network for container and premium last-mile deliveries. The platform manages the flow of goods from sea ports to warehouses, between distribution hubs, and directly to end-customers – supporting both industrial-scale freight and





Business Challenge

The client faced several critical issues in managing container and specialized deliveries across the U.S.:

To remain competitive, the client required a platform capable of dynamically sourcing available drivers and vehicles, minimizing idle time, and automating complex route and delivery decisions.

\checkmark

Lack of Real-Time Coordination:

The legacy systems lacked real-time insight into port schedules, traffic, or vehicle availability, which made proactive planning impossible and introduced bottlenecks across the delivery chain.

Driver and Truck Availability:

A recurring problem was the inability to quickly locate and assign available trucks and drivers for container pickup at sea ports. Containers often remained idle at terminals or port-side warehouses, incurring storage fees and causing downstream delays.

\checkmark

 \checkmark

Diverse Delivery Requirements:

The company needed a unified solution that could serve not only industrial-scale container logistics, but also handle white glove delivery scenarios — where precise timing, handling, and customer experience are essential.

\checkmark

Idle Container Costs:

Delays in finding transportation resulted in long dwell times for containers, leading to increased demurrage charges and inefficient asset utilization.





Technical Challenge

The core technical challenge was to build a logistics platform that could intelligently track, predict, and coordinate the availability of freight drivers effectively orchestrating deliveries in real time. This required solving several deep operational and engineering problems:



Dynamic Driver Availability Tracking:

Identify which drivers are currently active, who will be available later, and who might drop out — using real-time app signals and historical behavior patterns.



AI-Assisted Dispatcher Operations:

Provide dispatchers with AI assistance for interpreting freetext updates from drivers, generating responses, and flagging operational issues.



Disruption Response & Reallocation:

Automatically handle no-shows, early/late arrivals, and lastminute cancellations — and reassign tasks without human bottlenecks.



Document Processing & Validation:

Process and validate proof-of-delivery documents, customs paperwork, and gate passes via AI — with OCR and field extraction.



Model definitions

MODEL NAME	ARCHITECTURE	PROVIDER	PURPOSE	INPUT	OUTPUT	USAGE CONTEXT
DriverPredict-XGB	XGBoost (Gradient Boosting Trees)	Custom (in-house)	Predicts driver availability based on historical behavior and app data	Driver ID, activity history, shifts, app pings	Availability probability for time slots	Used in planning to pre-allocate or reassign deliveries
DispatchGPT	Large Language Model (GPT-4 Turbo)	OpenAl	Interprets driver messages, suggests replies, flags risk	Free-text message, dispatcher context	Suggested response, summarized intent, risk flags	Live assistant in dispatcher UI
DocOCR-Validator	TrOCR + LayoutLMv3	Microsoft / HuggingFace	Processes and validates delivery documents	Scans, PDFs, images	Structured data, validation results, compliance status	Document processing for proof-of-delivery and compliance workflows

Model Runtime Characteristics

MODEL NAME	AVG INFERENCE TIME	RETRAINING FREQUENCY	CONTEXT VOLUME	EXPLAINABILITY	REAL-TIME COMPATIBLE
DriverPredict-XGB	~200 ms	Weekly	~50 variables (driver activity patterns)	SHAP feature importance	Ves Yes
DispatchGPT	~500–700 ms	N/A (prompt-based)	~300–500 tokens (text + metadata)	Natural-language explanation (LLM)	Ves Yes
DocOCR-Validator	~2-3 sec	Rare (base models), rules updated quarterly	~1 document page (image + form structure)	A Partial (OCR confidence + heuristics)	A Batch or async preferred



Model Comparison for Dispatcher Assistant

MODEL	RESPONSE TIME	OUTPUT QUALITY	STRENGTHS	WEAKNESSES
GPT-4 Turbo (Chosen)	~500–700 ms	High (natural, contextual, actionable)	Fast, cost-effective, highly contextual	Limited to ~128k tokens; sometimes verbose
GPT-4 (original)	~2-3 sec	Very High (nuanced, detailed, reliable)	Deep understanding, great in edge cases	High latency and cost for frequent UI interactions

DriverPredict Model Comparison

MODEL NAME	ACCURACY	TRAINING TIME	INFERENCE TIME	EXPLAINABILITY	STRENGTHS	WEAKNESSES TIME
XGBoost (Chosen)	High (great for structured behavioral data)	Fast	~200 ms	SHAP feature importance	Interpretable, robust, effective for tabular data	May need manual feature engineering
LightGBM	Comparable to XGBoost	Very fast	~150 ms	SHAP	Fast training, scalable on large datasets	Slightly harder to tune for accuracy
Temporal Fusion Transformer	Very high for time- series	Slow	~1-2 sec	Attention-based insights	Captures time dependencies, interpretable via attention	Complex infra, slow inference, overkill





XGBoost



Express js





GPT-4 Turbo



Node.js

Tech Stack



TrOCR



LayoutLMv3



GraphQL



React



PostgreSQL

AWS (ECS, Fargate, Lambda)

aws



How Was Al Integrated?

The platform primarily operates as a logistics marketplace, where drivers — either independently or via affiliated companies — browse and accept freight assignments directly through the app. To ensure reliability, especially during peak load times or when certain shipments were not picked up, the company also maintained a pool of internal and contracted drivers who could fill the gaps.

Phase 1

Handling Driver Inquiries and Documents

The first use case for AI was to answer frequently asked questions from drivers, such as:

- "Where do I drop off the container?"
- "What documents are needed at pickup?"
- "Is the terminal open before 7am?"

These questions were common and predictable, though phrased in slightly different ways. Using GPT-4 Turbo, we trained an assistant to handle them based on historical driverdispatcher conversations. For document-related queries, the Al also leveraged OCR and document understanding models (TrOCR + LayoutLMv3) to extract key data from uploaded scans or images.

This freed up dispatchers from a high volume of routine messages and reduced driver response time significantly.

In earlier stages, when a shipment remained unassigned or was at risk of delay, dispatch managers would manually reach out to drivers, remind them to log in, or negotiate urgent coverage. While effective, this workflow placed a heavy burden on the operations team.

To reduce that load, AI was gradually introduced — starting with the most repetitive and scalable processes.



Phase 2 Predicting Assignment Gaps Before They Happen

Next, we focused on a more strategic application: helping managers detect in advance when a delivery might not be picked up through the marketplace.

A predictive model was trained to identify patterns in regions, load types, and time windows where assignments historically fell through. Instead of waiting until the last moment, the system could now surface shipments with a high likelihood of going unclaimed — even if the deadline hadn't passed.

By doing so, the system helped minimize disruptions and gave dispatchers more time and flexibility to act — turning reactive triage into proactive planning.

With AI embedded into day-to-day workflows, the platform now offloads a large portion of repetitive operational tasks:



Answers driver questions automatically via AI assistant



Interprets documents and extracts key delivery details

This shift allows managers to:



Focus on exceptions that require human judgment



Make more strategic decisions, not micromanage the routine

This enabled the platform to:



Send proactive reminders to drivers





Suggest fallback options from the internal driver pool

\bigcirc

Identifies high-risk assignments early based on patterns and rules



Trust the system to handle predictable, repeatable tasks



Our Solution

We delivered a modular, AI-driven logistics platform optimized for real-time container and white glove coordination. Key components include:

DriverPredict-XGB

A predictive model (XGBoost) that forecasts driver availability using real-time telemetry, app usage, and historical shift patterns.

Driver Name	Region	Vehicle	Status	Last Shift	Next Available	Availability	v Al Flag	Al Note
Ethan Carter	Queens, NY	Truck	Assigned	May 28, 2024, night	May 30, 2024	91%		Regular night shifts, low fatigue risk
E. Tanaka	Brooklyn, NY	Truck	Resting	May 27, 2024, day	May 31, 2024	0 77%	A Needs ping	Inactive for 2 days, may need confirmation before assignment
D. Patel	Long Island City, NY	Van	In Transit	May 28, 2024, night	June 2, 2024	89%		En route — ETA 1h; auto-update on arrival
M. Kovačević	Manhattan, NY	Box Truck	Available	May 30, 2024, day	June 1, 2024	0 85%		Reliable for time- sensitive premium routes
Lucas Grant	Williamsburg, NY	Van	Unassigned	May 30, 2024, morning	June 1, 2024	63%	A Reassignment	Awaiting reassignment; recently canceled task
E. Tanaka	East Harlem, NY	Box Truck	Assigned	May 28, 2024, day	May 30, 2024	94%		Consistently accepts morning industrial routes
Aiden Brooks	Bronx, NY	Box Truck	Resting	May 30, 2024, day	June 4, 2024	68%	🛦 Monitor load	Accepts short trips only; history of fatigue complaints
Carter James	Queens, NY	Van	Blocked	June 2, 2024, day	June 6, 2024	e 28%	High rejection	4 rejections in last 5 days; no-show flagged last week
Wyatt Cole	Staten Island, NY	Van	Assigned	May 30, 2024, day	May 31, 2024	87%		Strong acceptance rate; performs well in white glove deliveries

Driver Availability

DispatchGPT (GPT-4 Turbo)

An LLM-powered assistant that helps dispatchers interpret driver messages, suggest responses, and recommend reassignment actions.





Smart Reallocation Engine

Combines rule-based logic with reinforcement learning to dynamically reassign shipments based on real-time availability and delivery windows.



eason	Accepted By	ETA Impact
op timeout	Auto	+5 min
te check-in	Dispatcher	-2 min
hicle breakdown	Auto	TBD



DocOCR-Validator

A document-processing pipeline built on TrOCR + LayoutLMv3 for extracting and verifying fields in delivery documents and compliance forms.

Extracted Fields Panel

Field	Extracted Value	Validation Status	AI Confidence
Shipment ID	#SH984312	🥑 Valid	97%
Delivery Date	05/27/2025	🥝 Valid	94%
Receiver Signature	Not detected	🙁 Missing	
Client Name	A. Carter Logistics LLC	🔺 Needs Review	81%
Gate Pass Number	509842-A	🥑 Valid	92%
POD Document Type	Bill of Lading	🥑 Valid	88%
Delivery Address	2905 W 14th St, NY	🥑 Valid	90%
POD Time	03:45 PM, 05/28/2025	🥝 Valid	98%

Scalable Architecture

Backend (Node.js, GraphQL, AWS, PostgreSQL + MongoDB) and frontend (React) support large-scale dispatching and operator control.

Untimo			Ava Pochonco Timo		Poquests per Minut	
		act 20 days	292 mc (APLa)			
99.97%		ast ou days	JOZ MS (APIA)	vg, tast nour)	5,412 RPM	lavg tast nou
+0.01%			P95 latency: 680 ms		+100%	
Deserves	11	Thursday	0 220/ total	(lact hour)	Course and the	Chattan
Resource	Usage	Threshold	0.23% total	(last hour)	Component	Status
Resource CPU (avg)	Usage 61%	Threshold 80%	0.23% total API 5xx: AI fallback:	(last hour) 0.08% 0.05%	Component GraphQL Gateway	Status 🥑 Healthy
Resource CPU (avg) RAM (avg)	Usage 61% 72%	Threshold 80% 85%	0.23% total API 5xx: AI fallback: OCR Extraction Failures: Retries handled by queue buffer (auto	(last hour) 0.08% 0.05% 0.10% -recovery success: 96%)	Component GraphQL Gateway Eligibility Engine	Status 🥪 Healthy 🥑 Healthy
Resource CPU (avg) RAM (avg) Disk I/O	Usage 61% 72% Moderate	Threshold 80% 85% -	0.23% total API 5xx: AI fallback: OCR Extraction Failures: Retries handled by queue buffer (auto	(last hour) 0.08% 0.05% 0.10% -recovery success: 96%)	Component GraphQL Gateway Eligibility Engine Al Risk Validator	Status Healthy Healthy Slow
Resource CPU (avg) RAM (avg) Disk I/O GPU (Al node)	Usage 61% 72% Moderate 38% load	Threshold 80% 85% - 90%	0.23% total API 5xx: AI fallback: OCR Extraction Failures: Retries handled by queue buffer (auto	(last hour) 0.08% 0.05% 0.10% -recovery success: 96%)	Component GraphQL Gateway Eligibility Engine AI Risk Validator DocOCR Processor	Status Healthy Healthy Slow Healthy



Impact

Up to 95% automation

of route and assignment planning

99% on-time pickup

and delivery across **500,000+** annual shipments

3x improvement in SLA adherence

for white glove deliveries

70% boost in operational efficiency

for container transport

30% fewer idle/empty miles

reducing logistics overhead and CO₂ emissions

>90% automation

in document verification and compliance handling

