2026

# The Buyers Guide



# E-Commerce QA

A practical framework for evaluating agentic QA
platforms, running high-signal pilots, and building
the business case for faster, more reliable releases.

spur

# How to use this buyers guide

This guide provides a practical framework for evaluating modern QA approaches built for e-commerce. It focuses on how teams can assess their current process, define meaningful success metrics, and run high-signal pilots to determine whether a new class of testing platforms is a better fit for how their business operates today.

## What's Inside

# The Growing Gap in E-Commerce Testing

AI is making development faster than ever, but now the bottleneck has shifted to QA. E-commerce teams are impacted most, with storefronts changing daily through promotions, A/B experiments, and seasonal campaigns - yet most testing processes haven't kept up.

What we notice here is that teams are consistently struggling to automate testing fast enough to keep pace with AI-accelerated development

| 2 weeks | 80%+ | 10 FTEs |
|---|---|---|
| Avg. release cycle from dev complete to live | Regression testing effort that is still manual | Typical engineering time spent on QA |

→ *These estimates are based on conversations with brands doing approximately $60M+ in annual e-commerce revenue.*

## If the release process takes two weeks and still lets bugs reach production, the problem generally isn't the team—it's the tooling.

◄ **Back Market**

*"We would gather a team of engineers and manually 'walk the store' every Friday just to catch issues. Sometimes it would work, but most times a bug would slip through the cracks."*
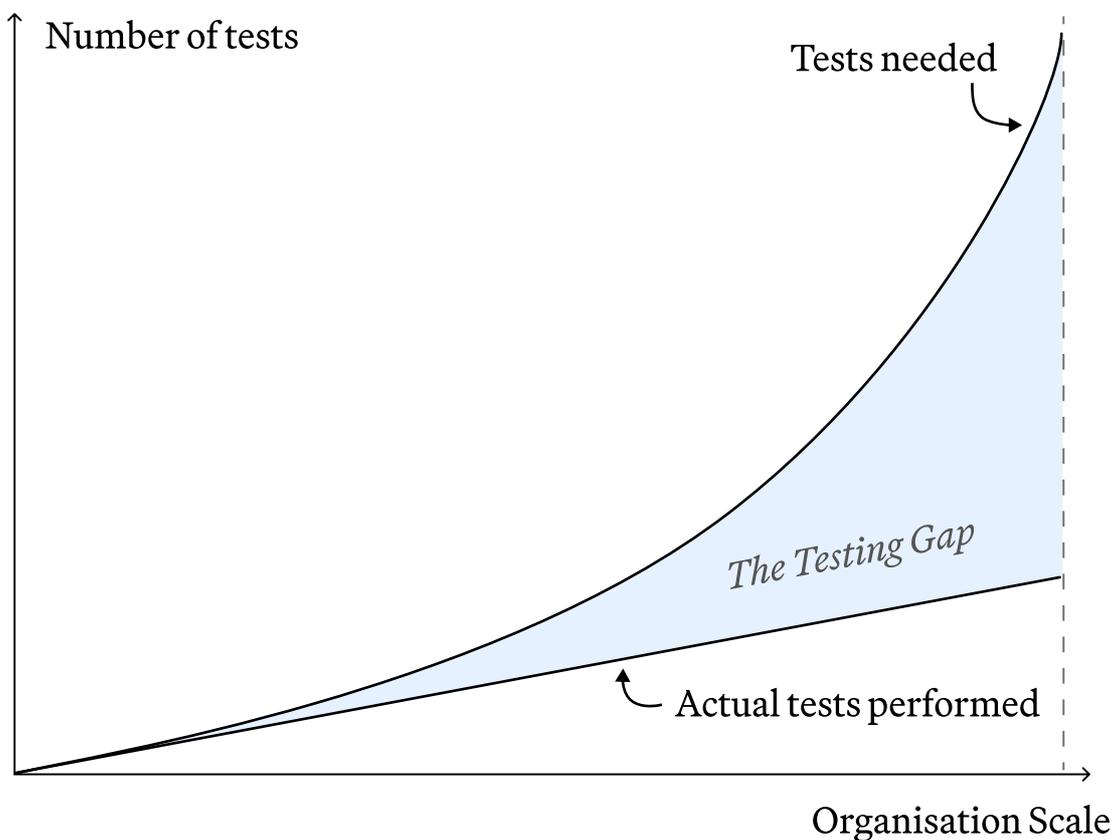
**Ray Ho**
VP of Product @ Back Market

# As teams scale, QA often becomes the bottleneck

As product surface area increases, the number of tests required to cover all screens and variations increases too. When teams reach the organizational maturity to start extensive A/B variant testing across multiple brands and locales, they often forego complex regression testing simply because the volume becomes unmanageable for human QA teams.



Unscalability of automated and manual testing

Now that AI is accelerating development speed, **the gap between how fast features ship and how well they're tested is widening at an alarming rate.** For large e-commerce brands juggling thousands of product pages, regional storefronts, and constantly changing checkout flows, that gap has real consequences: broken experiences, missed bugs, and compliance failures that only show up in production.

# Where current QA processes break down in E-Commerce

Most automation tools were built for **static, predictable web applications**. E-commerce storefronts are neither of those things.

A typical storefront has promotional banners that rotate weekly, product data that updates in real time, multiple layout variants running simultaneously through A/B experiments, and user flows that differ by geography, device, and customer segment.

That level of change is precisely what causes traditional scripted automation to fail.

| | |
|---|---|
| **Dynamic user interfaces** | Marketing teams push content updates, modal pop-ups, and seasonal campaigns on a weekly or even daily basis. Any test suites that relies on fixed element selector (Playwright, Cypress, Selenium) will break |
| **Volatile merchandising data** | Prices, inventory, and product details shift constantly.<br><br>A test that passed yesterday may flag a false failure today because a product went out of stock. |
| **A/B experiments and personalization** | The largest e-commerce brands often have more than 50 experiments simultaneously across their brands, changing depending on which test variant the visitor is in. |
| **Multi-locale and multi-brand complexity** | Global brands often operate dozens of storefronts with localized content, currencies, and checkout flows. Maintaining a separate test suite for each is expensive and difficult to scale. |
| **Visual and experiential quality** | A functional test may pass, but be riddled with UI issues such as content overflows, translation issues, and incorrect images. These directly impact conversion and brand perception. |

# QA doesn't need to be so hard.

# The Modern Approach to QA

Modern QA platforms don't just test. They understand intent.

Rather than relying exclusively on rigid scripts, they focus on validating intent and behavior exactly how a real shopper experiences the site and remain resilient to UI and content changes.

Tests are easier to author, more adaptable over time and better aligned with how e-commerce systems evolve.

# Teams using Agentic QA solutions are enjoying benefits

## Write tests faster

Describe what you want to test in natural language, with no need for custom code, selector chains, or framework expertise.

## Trust the results

Failures are tied to actual user-visible issues, not flaky selectors or timing glitches. When a test fails, it means something a customer would actually notice went wrong.

## Cover more

With lower authoring and maintenance costs, teams can afford to test user journeys they previously had to skip, including edge cases, non-English locales, and visual quality checks.
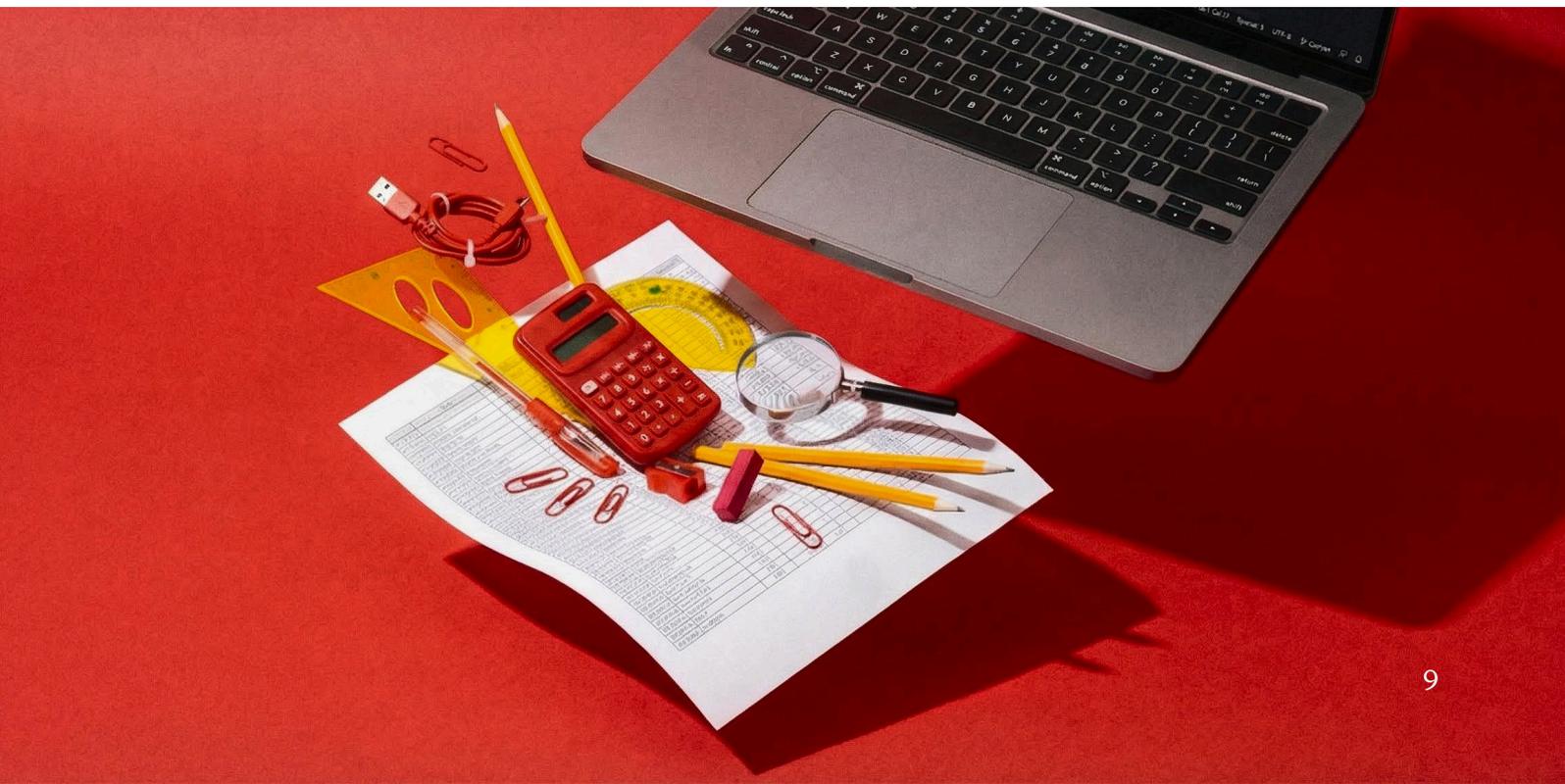
## Maintain less

Tests adapt to UI and content changes automatically. When marketing updates a banner or the product team reorganizes the checkout flow, existing tests keep working without manual intervention.

# Scripted automation testing vs. Agentic QA

This shift does not require replacing existing tools overnight. Most teams start with a focused pilot on their highest-impact flows and expand from there.

| Dimension | Agentic QA | Scripted Testing |
| --- | --- | --- |
| Time to author | Minutes to hours per test | Days to weeks per test |
| Maintenance | Minimal; self-heals automatically | Constant; breaks on UI changes |
| Stability | Consistent across runs | Flaky under change |
| Team accessibility | Anyone can author tests | Requires automation engineers |
| Visual validation | Built-in experience checks | Not supported natively |
| Dynamic content | Adapts to what's on screen | Breaks on variant changes |

# Where does your testing process stand?

Before evaluating new platforms, take an honest look at where your current process delivers and where it falls short. The goal is not to find fault, but rather to identify where better tooling would have the highest impact.

**Rate each statement from 1 (no confidence) to 5 (fully confident):**

| Criteria | Score (0 – 5) |
|---|---|
| **Release velocity**<br>Releases ship within 48 hours of dev-complete. Testing does not routinely add days to every sprint. | |
| **Regression burden**<br>Regression testing is primarily automated, not manually repeated each cycle. | |
| **Automation stability**<br>Automation stays stable through UI and content changes without constant maintenance. | |
| **Process maturity**<br>Test failures clearly indicate what broke and why, with a clear triage workflow, not a spreadsheet no one revisits. | |
| **Coverage**<br>Critical checkout and post-purchase flows have full coverage | |
| **Production escapes**<br>Bugs rarely escape to production. Functional errors, visual regressions, and checkout issues are caught before release. | |
| **Scalability**<br>QA scales across locales and devices without duplicating effort. | |

Scores consistently below 3 indicate meaningful room for improvement.

# Not every QA platform is built with e-commerce in mind.

Choosing the right QA partner is important.

General-purpose automation tools may work for testing a login form, but they tend to fall short when faced with the pace and  complexity of a modern storefront.

When evaluating solutions, look beyond marketing claims and focus on whether the platform actually handles the things that make e-commerce testing hard.

# Success Metrics: Setting the right benchmarks

Before kicking off a pilot, align on what a good fit for your team looks like. Defining metrics upfront creates a shared language across QA, engineering, product, and leadership. This prevents the evaluation from drifting into subjective impressions.

| Metric | Why it matters | ☑ Suggested Targets |
|---|---|---|
| Test authoring time | Faster authoring reduces adoption friction and accelerates coverage | Simple ≤ 15 min<br>Complex ≤ 45 min |
| Maintenance effort | Separates tools that save time from those that shift the work. | ≤ 10 min/test |
| Self Healing | Separates tools that save time from those that shift the work. | ≥ 60% auto-resolved |
| False positive rate | High false-positive rates erode trust and train teams to ignore results. | ≤ 2% per 100 runs |
| Run stability | Flaky tests create noise that masks real problems. | ≥ 98% pass consistency |
| Defect detection | Measures whether the tool catches what actually matters to customers. | ≥ 90% of seeded defects |
| Coverage breadth | Core journeys, localized flows, device-specific paths, post-purchase. | 100% of scoped flows |
| Time savings | The metric that builds the ROI case for leadership. | ≥ 40% reduction in QA cycle time |

These metrics translate testing performance into language the business understands. When the data shows that a new tool reduce authoring time by 80% and eliminated a week of manual regression per release, the path forward to testing for your team becomes clear.

# Running a Pilot That Gives Real Answers

A well-designed pilot or proof of concept is the fastest path to making the right decision for your brand. The goal is not exhaustive testing over months but to generate enough signal, in a short enough window, with which your team can move forward with conviction.

Here's what you should should aim to do in your pilot:

## Keep it focused.

One to two weeks, concentrated on your most critical flows. Plans and pricing, checkout, upgrades, cancellations. Test the flows where bugs hurt the most.

## Use production-like environments.

Staging with synthetic data reveals very little. Test against environments that mirror what your customers actually see.
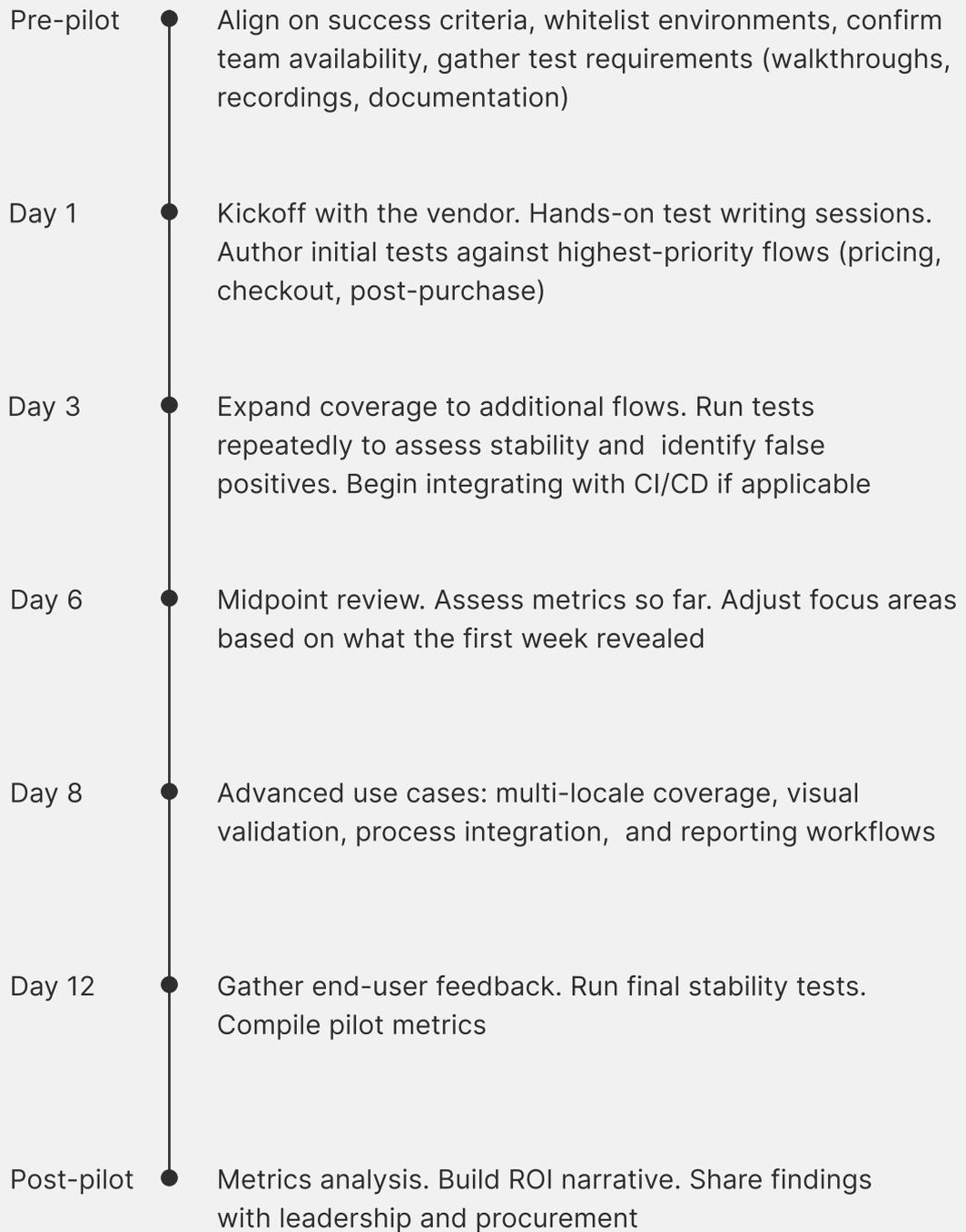
## Dedicate a team.

Assign 2–3 members who will work with the platform daily. Background pilots with sporadic check-ins rarely produce clear signal.

## Track metrics from day one.

Collect data on authoring time, maintenance, false positives, and defect detection throughout the pilot, not just at the end.

# A well designed pilot

**Pre-pilot** — Align on success criteria, whitelist environments, confirm team availability, gather test requirements (walkthroughs, recordings, documentation)

**Day 1** — Kickoff with the vendor. Hands-on test writing sessions. Author initial tests against highest-priority flows (pricing, checkout, post-purchase)

**Day 3** — Expand coverage to additional flows. Run tests repeatedly to assess stability and identify false positives. Begin integrating with CI/CD if applicable

**Day 6** — Midpoint review. Assess metrics so far. Adjust focus areas based on what the first week revealed

**Day 8** — Advanced use cases: multi-locale coverage, visual validation, process integration, and reporting workflows

**Day 12** — Gather end-user feedback. Run final stability tests. Compile pilot metrics

**Post-pilot** — Metrics analysis. Build ROI narrative. Share findings with leadership and procurement

→ Short, intensive pilots surface real performance data faster than extended evaluations with sporadic engagement. Two focused weeks will always beat two scattered months.

# Making the case for change

Once the pilot is completed, your decision shifts from evaluating the tool to now making the case for adopting it. The most successful teams treat this phase with the same rigor as the pilot itself.

Before the pilot began, the team should have aligned on success criteria and target metrics (see "Setting the Right Benchmarks"). The following steps turn that baseline into a decision.

Download your template ↗

# Making the case for change

Once the pilot is completed, your decision shifts from evaluating the tool to now making the case for adopting it. The most successful teams treat this phase with the same rigor as the pilot itself.

Before the pilot began, the team should have aligned on success criteria and target metrics (see "Setting the Right Benchmarks"). The following steps turn that baseline into a decision.

Immediately after the pilot:

## 1. Complete the success metrics scorecard

Return to the benchmarks defined before the pilot and fill in actual results while the data is fresh.

| Metric | Suggested Targets | Actual |
|---|---|---|
| Test authoring time | Simple ≤ 15 min<br>Complex ≤ 45 min | |
| Maintenance effort | ≤ 10 min/test | |
| Self Healing | ≥ 60% auto-resolved | |
| False positive rate | ≤ 2% per 100 runs | |
| Run stability | ≥ 98% pass consistency | |
| Defect detection | ≥ 90% of seeded defects | |
| Coverage breadth | 100% of scoped flows | |
| Time savings | ≥ 40% reduction in QA cycle time | |

## 2. Build the ROI narrative

Rate each statement from 1 (not at all) to 5 (strongly agree):

| Statement | Score |
| --- | --- |
| The platform delivered meaningful coverage within the pilot window | |
| Maintenance effort was materially lower than with existing tools | |
| Failures were actionable and tied to real customer-facing issues | |
| The team can see this platform as part of their daily workflow | |
| The business case justifies the investment | |

## 3. Document qualitative feedback

Honest feedback carries weight with stakeholders who don't read dashboards. Capture responses from the pilot team:

What was the team's first impression of the platform?

How did debugging compare to existing tools?

Which flows were easiest to author? Which were hardest?

Did the platform fit into the daily workflow or feel like extra work?

What would the team change about the pilot experience?

Would the team want to use this tool going forward? Why or why not?

# How Spur performs against these metrics

The frameworks discussed in this guide are designed to be vendor-neutral. Based on our current customers, here is how Spur performs when compared against the success metrics.

# Spur exceeds every benchmark that matters

### Test Authoring

## 2-10 mins
for simple flows

Plain-language, Video/Loom, and AI auto test creation. Complex flows under 30 mins, with no scripting expertise required.

Target: 20mins

### Maintenance

## Near 0%
maintenance effort

Agentic AI adapts to UI and content changes without manual updates. Tests survive redesigns.

Target: < 10min/test

### Self-Healing

## 90%+
auto-resolved

The agent interprets the screen like a customer would. Layout changes, modals, and content updates are handled automatically.

Target: < 10min/test

### False positives

## < 2%
per 100 runs

Failures tied to real user-visible issues, not broken selectors or timing problems

Target: ≤ 2% per 100 runs

### Run stability

## >99%+
pass consistency

Consistent results across repeated executions. No flaky runs undermining CI/CD confidence.

Target: ≥ 98%

### Defect detection

## 95%+
of seeded defects catch rate

Catches checkout errors, pricing issues, and post-purchase failures with screenshots, recordings, and reproduction paths.

Target: ≥ 90%

### Coverage

## 100%
of scoped flows

One test definition scales across devices, browsers, locales, and brands.

Target: ≥ 80%

### Time savings

## 60–80%
reduction in QA cycle time

Brands starting at 20–30% automation routinely double or triple coverage within weeks.

Target: ≥ 40% reduction

### QA Cycle

## 5-50×
faster release cycles

Faster release cycles with comprehensive test coverage from day one.

# The fastest way to evaluate these claims is to see Spur run on your ecommerce website and mobile app.

Our customers are already experiencing the ROI by having QA testing on autopilot.

**alo**

24 h → 2 h

Time taken for regression turnaround

**HELLO FRESH**

0 → 95%

Test coverage achieved in 2 weeks (production, web & mobile)

"Spur helped us build up our QA program by cutting manual QA time down from days to 30 mins each release. It's one of the major AI wins in our company, we went from 0 to 80% coverage in 1 month."

**Pete Franco**
President, LivingSpaces

**LIVING SPACES**

See more results like these in our case studies →

# Ready to see what Agentic QA looks like for your team?

Book a demo and we'll run Spur on your live site. You will see real test results on your workflows before making any commitment.

**Book a demo** ⧉