

# ISTANBUL SOFTWARE TESTING CONFERENCE

## Lessons Learned from AI-Powered Visual Reasoning

Detecting UI issues, accessibility problems,  
and design inconsistencies with LLMs

Risko Ruus | ISTC 2026

#ISTC2026



## 2006: Way Back When

### Automated Visual Testing Symbian Operating System

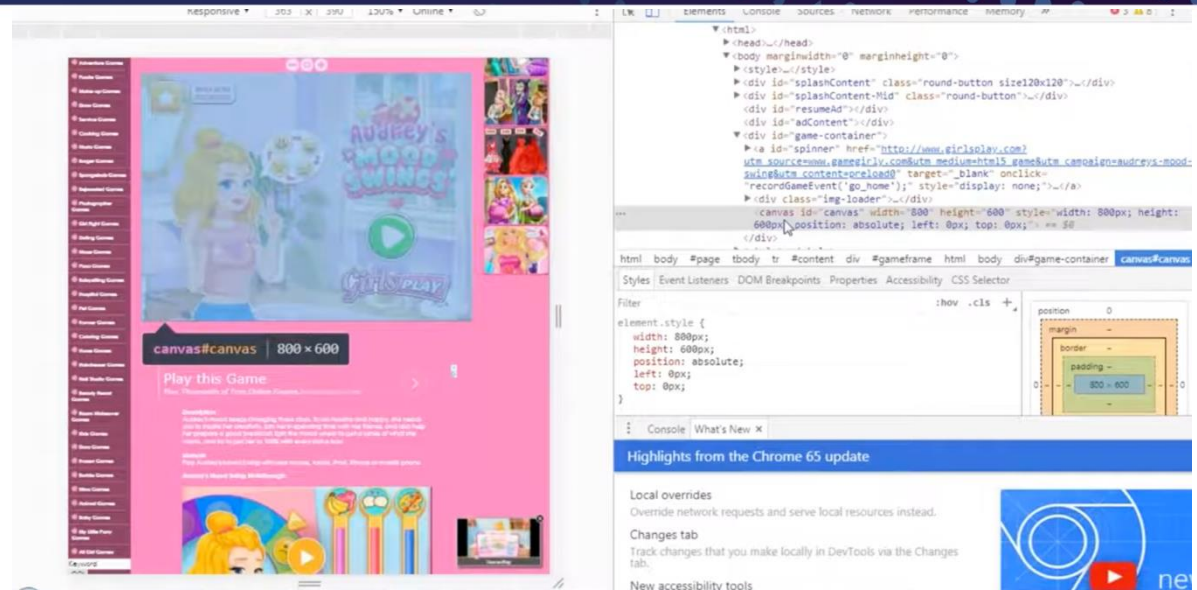
- Minor pixel variations cause assertion failures
- Unique images required for every device
- High maintenance demands



# 2010: Adobe Flash & Sikuli

## Automating Adobe Flash Content with Sikuli (OpenCV)

- Image recognition to identify GUI components
- Flash objects lacked IDs for UI automation
- Integrated it with Selenium WebDriver
- Image maintenance was a burden



# 2026: Similar Problems Still Today

---

## The Same Problem, New Tech

### Web: WebAssembly + Canvas

- App logic runs as opaque WASM bytecode — no DOM, no inspectable state
- UI rendered to a <canvas> — Selenium/Playwright see one element, not buttons or text

### Native: Flutter, Rive, game engines, etc.

- App paints its own surface via Skia/Metal/GL — no native view hierarchy
- Appium/Espresso/XCUITest see one container view — widgets inside are invisible

*Different stack, same testability gap: no selectors, pixel-based interaction, visual-only validation.*

# About Me

---

## Risko Ruus

Principal QA Engineer  
@Rush Street Interactive

## EXPERIENCE

20+ years in QA — from early-stage startups to products at scale

## HOBBIES

- Backcountry hiking
- Learning Japanese with my son
- Tinkering with AI



# Agenda

01 Traditional UI Assertions

02 Visual AI Reasoning  
Core Concepts &  
Basic Usage

03 Advanced Analysis  
Techniques

04 Integrating Visual  
Reasoning Checks into  
Workflows

01

---

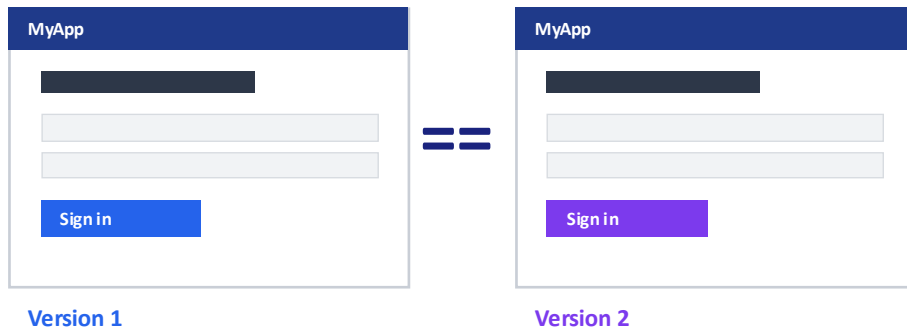
# Traditional UI Assertions

Traditional UI Assertion Approaches | Limitations

# Problems with Traditional Visual Check Approaches

---

## Approach 1: Full Match



### Problems:

- ⚠ Too flaky — fails on minor changes
- ⚠ Font rendering, anti-aliasing cause false failures
- ⚠ Requires constant screenshot updates

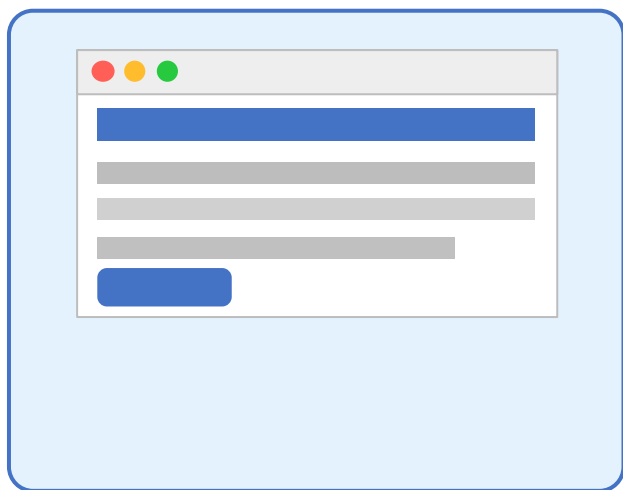
High Maintenance

Brittle

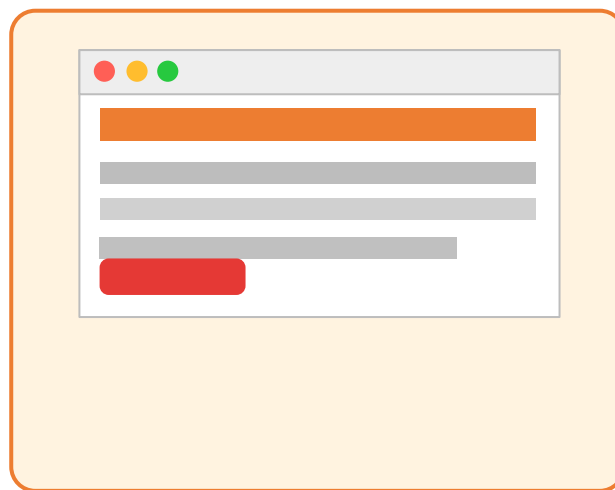


# How Traditional Pixel Diff Tools Work

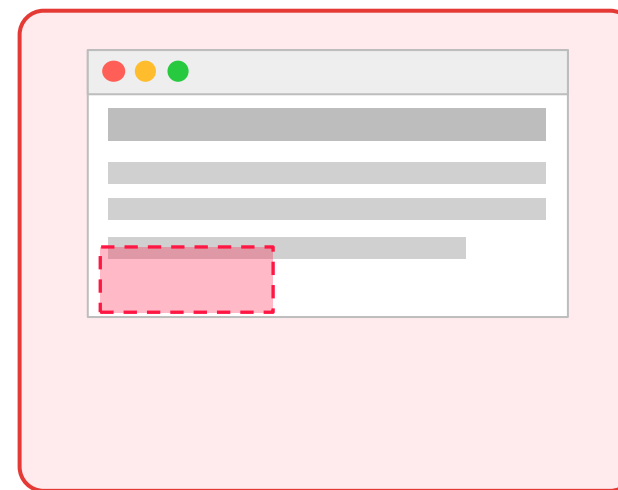
Applitools Eyes, Percy, and similar tools



Baseline Screenshot



Current Screenshot



Pixel Diff Output

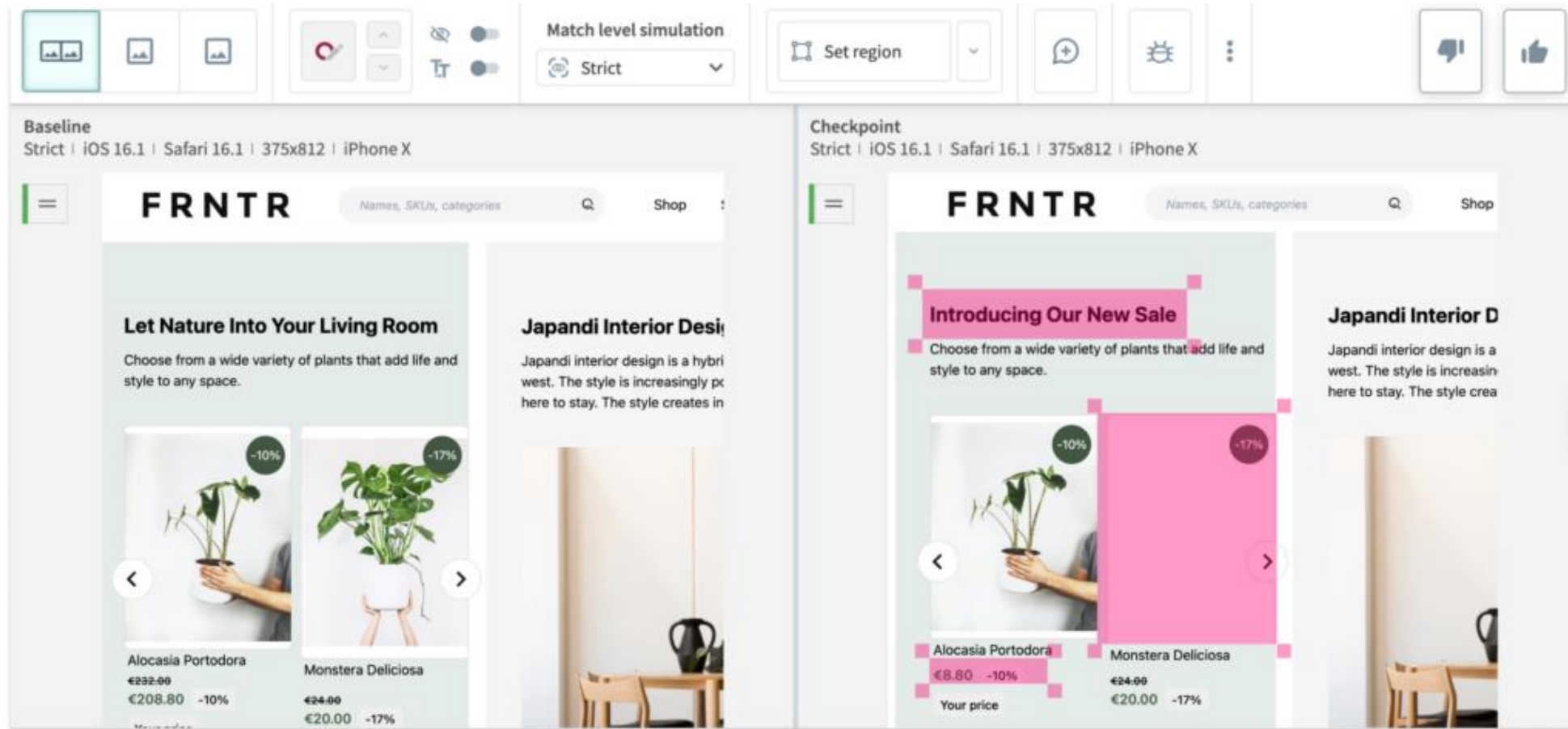
## What pixel diff reports:

- "12 pixels changed at region (140, 280)"
- "0.3% visual difference detected"
- "Button color shifted from #44C to #E53"

## What it can't tell you:

- ✗ Is the button still usable?
- ✗ Does the new color meet contrast standards?
- ✗ Was this change intentional or a bug?

# Example Pixel Diff Output



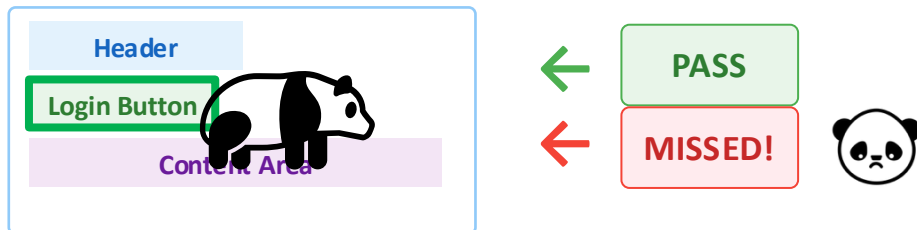
Applitoools Eyes Image Comparison

# Problems with Traditional Visual Check Approaches

## The Giant Panda Problem

### Approach 2: Specific Element Check

E.g. verify that a login button is displayed



### Problems:

- ⚠ Explicit checks miss new erroneous content
- ⚠ Only validates what you coded it to check

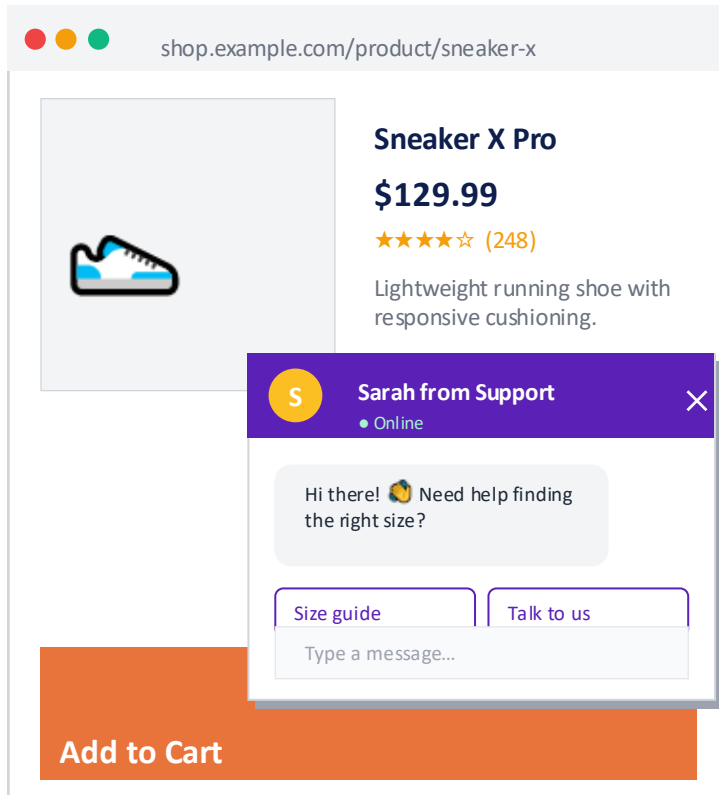
Blind Spots

False Confidence



# Real-World Examples: When isDisplayed/isVisible Lies

Each page looks fine to a non-vision UI test— but a real user has a poor user experience



shop.example.com/product/sneaker-x

**Sneaker X Pro**  
**\$129.99**  
★★★★☆ (248)  
Lightweight running shoe with responsive cushioning.

**Sarah from Support**  
● Online

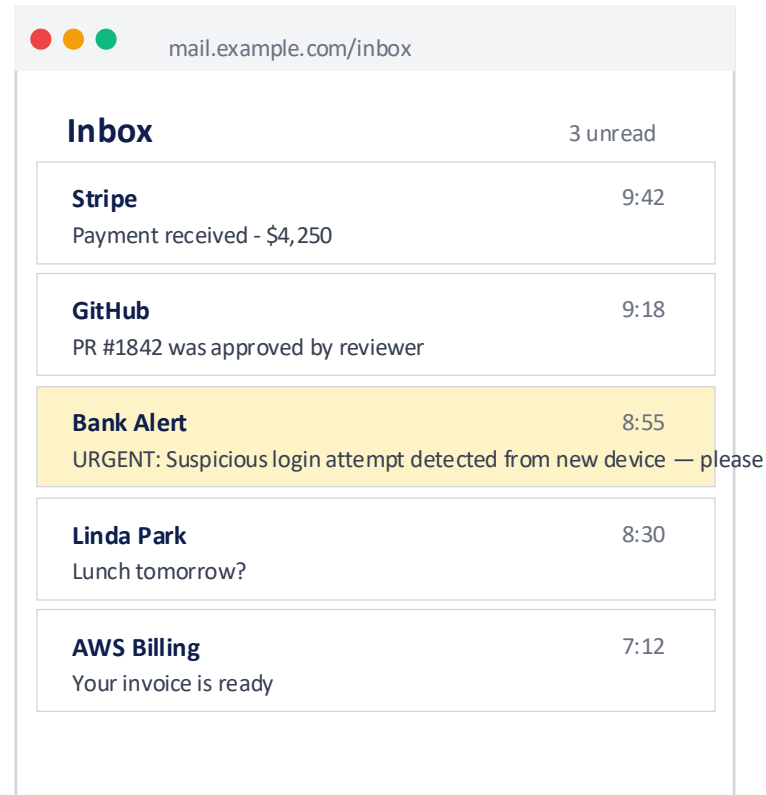
Hi there! 🧐 Need help finding the right size?

[Size guide](#) [Talk to us](#)

Type a message...

**Add to Cart**

Covered by overlay



mail.example.com/inbox

**Inbox** 3 unread

**Stripe** 9:42  
Payment received - \$4,250

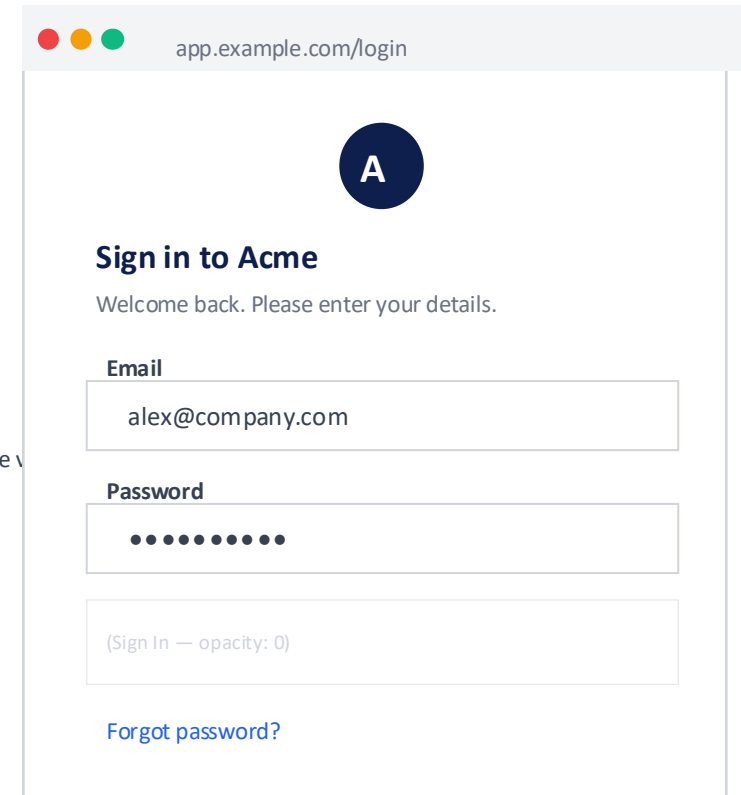
**GitHub** 9:18  
PR #1842 was approved by reviewer

**Bank Alert** 8:55  
URGENT: Suspicious login attempt detected from new device — please v

**Linda Park** 8:30  
Lunch tomorrow?

**AWS Billing** 7:12  
Your invoice is ready

Clipped by overflow



app.example.com/login

**A**

**Sign in to Acme**  
Welcome back. Please enter your details.

**Email**  
alex@company.com

**Password**  
●●●●●●●●

(Sign In — opacity: 0)

[Forgot password?](#)

Opacity: 0

02

---

# Visual AI Reasoning Core Concepts & Basic Usage

Traditional vs AI-powered testing | Analysis types

# How the AI Reviews a Screenshot

Three steps from screen image to readable feedback on bugs and design

## 1. Screenshot + prompt in

### Just a screenshot + text

Show the AI a picture with a task of finding issues. No design spec or "before" image needed, but it can be helpful.

## 2. Visual + semantic check

### Looks and understands

The AI reads the image alongside your instructions, spotting visual bugs (overlaps, cut-off text, contrast) and design issues using what it learned during training.

## 3. Clear feedback out

### A list of issues

Returns issues with severity and a short explanation. References shapes or regions by name, not pixel coordinates.

**Key insight:** *General-purpose LLMs catch many visual and design bugs out of the box — no fine-tuning required.*

# Example Vision LLM Input Image

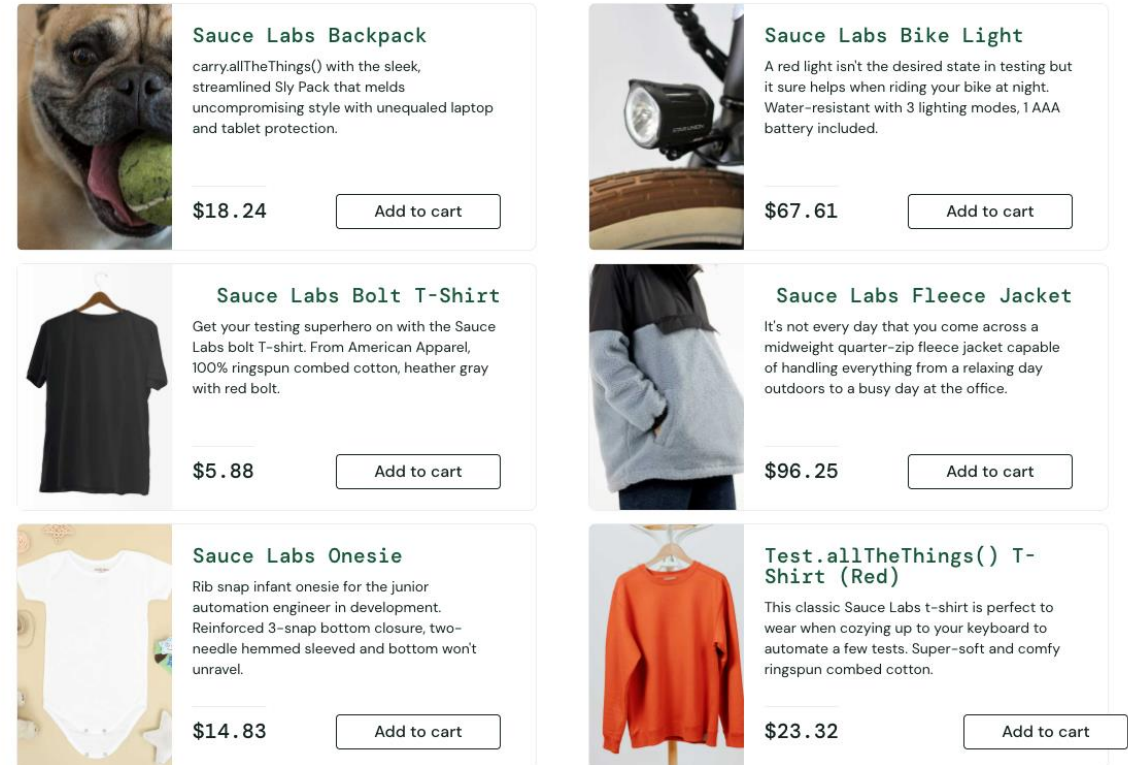
This public website was developed by Sauce Labs for demo/testing purposes:

<https://www.saucedemo.com/inventory.html>

It includes deliberate bugs that we expect the AI to identify.

**NB! Always check which data you are allowed to send to LLM providers e.g. enterprise / regulated / private data usage**

Swag Labs



Product Name	Description	Price	Action
Sauce Labs Backpack	carry.allTheThings() with the sleek, streamlined Sly Pack that melds uncompromising style with unequalled laptop and tablet protection.	\$18.24	Add to cart
Sauce Labs Bike Light	A red light isn't the desired state in testing but it sure helps when riding your bike at night. Water-resistant with 3 lighting modes, 1 AAA battery included.	\$67.61	Add to cart
Sauce Labs Bolt T-Shirt	Get your testing superhero on with the Sauce Labs bolt T-shirt. From American Apparel, 100% ringspun combed cotton, heather gray with red bolt.	\$5.88	Add to cart
Sauce Labs Fleece Jacket	It's not every day that you come across a midweight quarter-zip fleece jacket capable of handling everything from a relaxing day outdoors to a busy day at the office.	\$96.25	Add to cart
Sauce Labs Onesie	Rib snap infant onesie for the junior automation engineer in development. Reinforced 3-snap bottom closure, two-needle hemmed sleeved and bottom won't unravel.	\$14.83	Add to cart
Test.allTheThings() T-Shirt (Red)	This classic Sauce Labs t-shirt is perfect to wear when cozying up to your keyboard to automate a few tests. Super-soft and comfy ringspun combed cotton.	\$23.32	Add to cart

# Example Vision LLM Prompt

## SYSTEM PROMPT

You are a senior QA engineer. Analyze this e-commerce screenshot and identify all UI, content, and functional issues.

### Severity levels:

**Critical** – Blocks purchasing or damages trust

**Major** – Clear defect, workaround exists

**Minor** – Cosmetic, unlikely to affect users

### Check for:

Broken/placeholder content, layout issues, wrong images, suspicious pricing, accessibility, staging data leaks.

*Output valid JSON using this schema →*



## JSON SCHEMA

```
{
  "issues": [
    {
      "id": 1,
      "severity": "critical|major|minor",
      "category": "Content, UI/Layout...",
      "element": "Specific element",
      "description": "What is wrong",
      "expected": "Correct state",
      "impact": "User/biz impact"
    }
  ],
  "summary": {
    "critical": 0,
    "major": 0,
    "minor": 0,
    "total": 0,
    "overall_assessment": "Summary..."
  }
}
```

# Example Vision LLM Output

JSON OUTPUT

```
{  
  "id": 1,  
  "severity": "critical",  
  "category": "Content/Copy",  
  "element": "Page title and all product names",  
  "description": "The entire storefront uses 'Sauce Labs' branding and product names (e.g., 'Sauce Labs Backpack', 'Sauce Labs Bolt T-Shirt'). Sauce Labs is a test automation platform; this is their demo app leaking test data into what should be production.",  
  "expected": "Real brand name and real product names for an actual production store",  
  "impact": "Destroys customer trust entirely; no real customer would purchase from a store obviously running test data. Indicates staging or QA environment exposed publicly."  
},
```


# Example Vision LLM Output

JSON OUTPUT

```
{
  "id": 2,
  "severity": "critical",
  "category": "Content/Copy",
  "element": "Product name: 'Test.allTheThings() T-Shirt (Red)'",
  "description": "Product name contains a literal test method call 'Test.allTheThings()' – unmistakable test/developer placeholder data visible to end users.",
  "expected": "A real, customer-facing product name (e.g., 'Classic Red T-Shirt')",
  "impact": "Confirms test data is in production; severely damages credibility and trust. Indicates staging or QA environment exposed."
}
```

# Agentic Vision with Annotations Highlighting Detected Issues (Gemini 3 Flash)


Swag Labs



**1** Sauce Labs Backpack

carry.allTheThings() with the sleek, streamlined Sly Pack that melds uncompromising style with unequaled laptop and **5** t protection.


\$18.24



Sauce Labs Bike Light

A red light isn't the desired state in testing but it sure helps when riding your bike at night. Water-resistant with 3 lighting modes, 1 AAA battery included.


\$67.61



**2** Sauce Labs Bolt T-Shirt

Get your testing superhero on with the Sauce Labs bolt T-shirt. From American Apparel, 100% ringspun combed cotton, heather gray with red bolt.


\$5.88



**6** Sauce Labs Fleece Jacket

It's not every day that you come across a **6** midweight quarter-zip fleece jacket capable of handling everything from a relaxing day outdoors to a busy day at the office.


\$96.25



Sauce Labs Onesie

Rib snap infant onesie for the junior automation engineer in development. Reinforced 3-snap bottom closure, two-needle hemmed sleeved and bottom won't unravel.

\$14.83



**3** Test.allTheThings() T-Shirt (Red)

This classic Sauce Labs t-shirt is perfect to wear when cozying up to your keyboard to automate a few tests. Super-soft and comfy ringspun combed cotton.

\$23.32

**4**

ID	Severity	Issue Summary
1	Critical	Wrong image on Backpack (shows dog photo)
2	Major	Bolt T-Shirt missing red bolt graphic
3	Major	T-Shirt shows sweatshirt image
4	Minor	Inconsistent Add to Cart button width
5	Minor	Developer placeholder text in description
6	Major	Suspicious pricing suggests staging data

# From Pixel Diffs to LLM Visual Intelligence Comparison



## Traditional Pixel Diff

vs

## AI-Powered Visual Reasoning

<b>Detects</b>	Pixel-level differences	Semantic UI issues
<b>Reports</b>	Changed vs unchanged pixels	Layout, hierarchy, meaning, accessibility
<b>False Positives</b>	High (any render diff over a threshold)	Medium (contextual reasoning)
<b>Feedback</b>	"Pixels differ at (x,y)"	"Button is hidden by overlay"
<b>Exclusions</b>	Manual exclusion areas per image	Ignore prompts
<b>Speed</b>	Milliseconds (local compute)	Seconds per image (depends on the model reasoning level)
<b>Cost</b>	Near zero (local CPU)	API token fees per invocation
<b>Determinism</b>	Deterministic	<b>Non-deterministic</b>

AI visual checks should complement, not replace, deterministic checks!

03

---

# Advanced Analysis Techniques

Model selection | Agent review | Reducing noise and improving accuracy

# Model Overview – Current Popular Vision-Capable LLMs\*

## OpenAI

GPT-5.5  
GPT-5.4 Pro  
GPT-5.4  
GPT-5.4 Mini  
GPT-5.4 Nano

## Anthropic

Opus 4.7  
Sonnet 4.6  
Haiku 4.5

## Google

Gemini 3.5 Flash  
Gemini 3.1 Pro  
Gemini 3.1 Flash-Lite  
Gemini 3 Flash

## Open Weighted Models

Qwen 3.5 397B  
A17B  
Kimi K2.6  
Qwen3.5 27B  
...

### Thinking effort is configurable on most of these models

A "thinking effort" or "reasoning budget" setting (e.g. minimal / low / medium / high, or an explicit token cap).

It controls how many internal reasoning tokens the model is allowed to spend before answering.

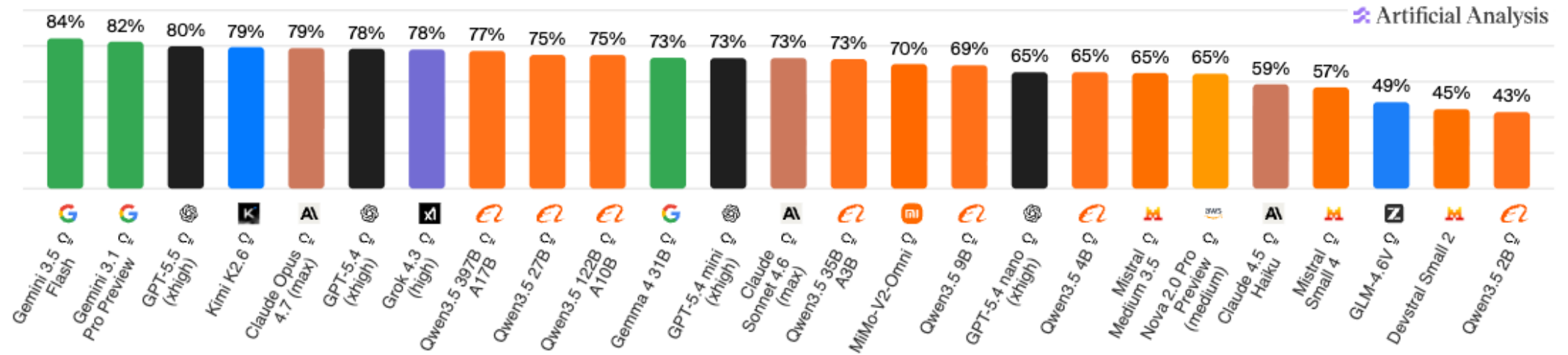
Higher effort = better accuracy on hard tasks, at the cost of more latency and tokens; lower effort = faster, cheaper responses for simple queries.

# Visual Reasoning Intelligence (MMMU Pro evaluation)

25 of 190 models ⌵ ⌵ ↺ 🗖

0 of 417 model & provider com... ⌵

Visual reasoning intelligence: MMMU Pro evaluation



💡 Reasoning models are indicated by a lightbulb icon.

① **MMMU Pro** ✕

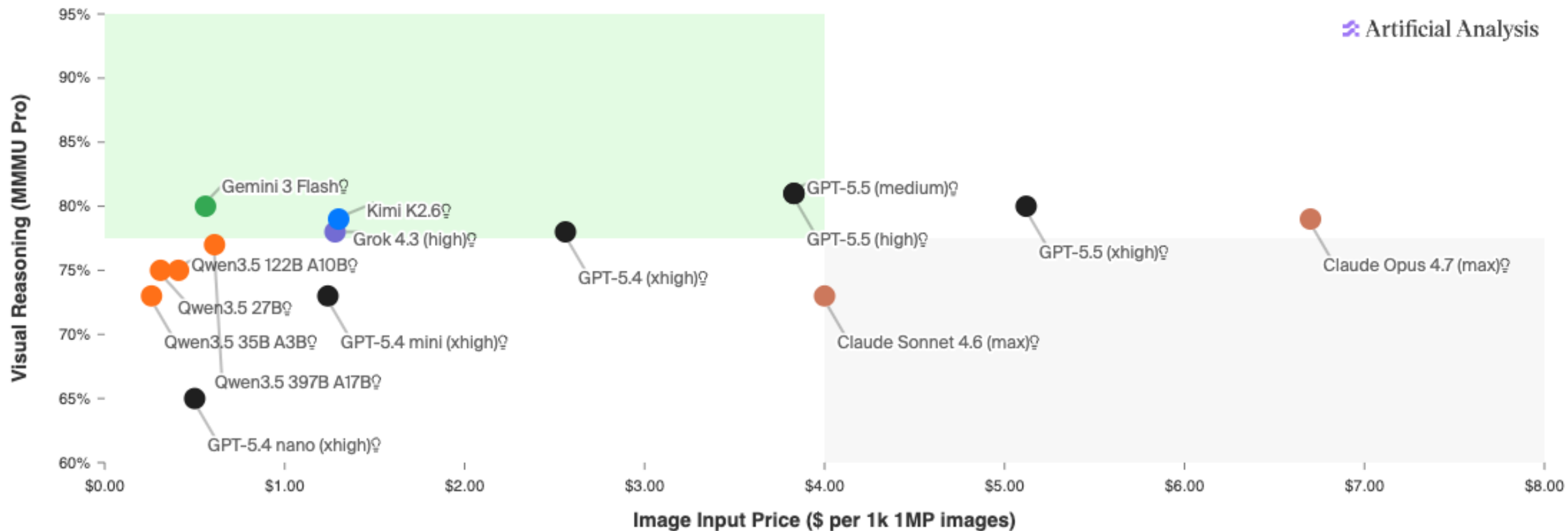
Multimodal reasoning quality evaluation based on 1.7k questions which require interpreting and reasoning over images.

# Vision Models Comparison – Visual Reasoning vs. Price\*

Visual reasoning intelligence: MMMU Pro evaluation - Image input price: USD per 1k images at 1MP (1024x1024)

Most attractive quadrant

- Claude Opus 4.7 (max)
- Claude Sonnet 4.6 (max)
- Gemini 3 Flash
- GPT-5.4 (xhigh)
- GPT-5.4 mini (xhigh)
- GPT-5.4 nano (xhigh)
- GPT-5.5 (high)
- GPT-5.5 (medium)
- GPT-5.5 (xhigh)
- Grok 4.3 (high)
- Kimi K2.6
- Qwen3.5 122B A10B
- Qwen3.5 27B
- Qwen3.5 35B A3B
- Qwen3.5 397B A17B



\*as of: May 20, 2026

<https://artificialanalysis.ai/models/multimodal/vision>

# Vision Models Comparison – Visual Reasoning vs. Latency\*

Visual reasoning intelligence: MMMU Pro evaluation - Latency (time to first token)

Most attractive quadrant

- Claude Opus 4.7 (max)
 ■ Claude Sonnet 4.6 (max)
 ■ Gemini 3 Flash
 ■ Gemini 3.1 Pro Preview
 ■ Gemini 3.5 Flash
 ■ Gemma 4 31B
 ■ GPT-5.4 (xhigh)
- GPT-5.4 mini (xhigh)
 ■ GPT-5.4 nano (xhigh)
 ■ GPT-5.5 (high)
 ■ GPT-5.5 (medium)
 ■ GPT-5.5 (xhigh)
 ■ Grok 4.3 (high)
 ■ Kimi K2.6
- Qwen3.5 122B A10B
 ■ Qwen3.5 27B
 ■ Qwen3.5 35B A3B
 ■ Qwen3.5 397B A17B

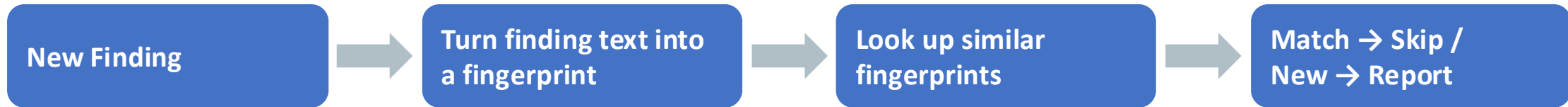


\*as of: May 20, 2026

<https://artificialanalysis.ai/models/multimodal/vision>

# Building a findings database – Dealing with duplicates

---



## Why Deduplicate?

---

- LLMs re-flag the same issue with different wording across runs
- Reviewers waste time closing duplicates instead of fixing bugs
- Plain text matching misses rephrased findings
- We need to compare by meaning, not just words

## How It Works

---

### Summarize each finding

Convert the issue into a numeric "meaning fingerprint"

### Save it to a database

Store fingerprints alongside each issue for fast lookup

### Compare new findings

If the new fingerprint closely matches an existing one, treat it as a duplicate and link to the original.

# Example: Same Bug, Different Words

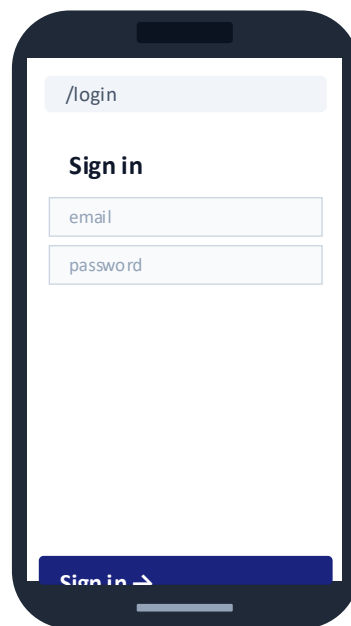
The AI inspects the app twice and describes the same problem in two ways. An embedding model recognizes they're the same and groups them together.

## Finding A — Run 1 (build #4821)

*"Login button on mobile is cut off at the bottom; users can't tap submit."*

## Finding B — Run 2 (build #4889)

*"Sign-in CTA overflows viewport on iPhone — primary action not clickable."*



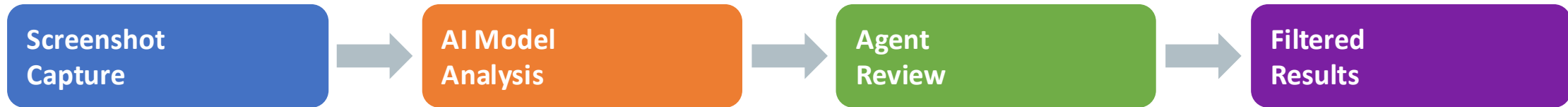
94%  
match



*The second finding is similar enough to be linked to the existing group — no duplicate ticket created.*

# Agent Review & False Positive Filtering

---



## How Agent Review Works

---

- An AI Agent or multiple agents review the findings
- Filters out noise and low-confidence detections
- Validates issues against domain-specific rules
- Produces a curated, actionable report

## When to Enable / Disable

---

### Enable when:

- High false positive rate in initial results
- Domain-specific rules need enforcement

### Disable when:

- Speed is critical (adds ~ an additional check)
- Cost becomes too expensive
- Initial model already has high enough accuracy

04

---

# Integrating Visual Reasoning Checks into Workflows

Test automation & CI/CD integration

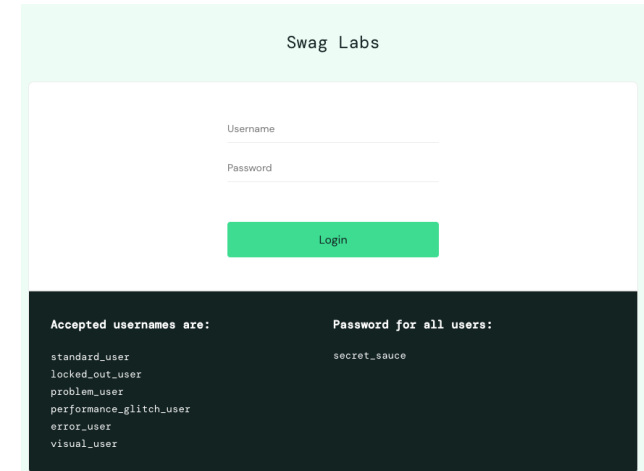
# Integrating Visual Reasoning Checks

Selenium / Playwright / Cypress / Appium

Add as a library for existing test projects

```
test_visual_check.e2e.ts
```

```
const result = await ai.check(screenshot, [  
  "The login button is visible",  
  "No error messages are displayed",  
  "A banner ad is present at the top of the page"  
]);  
  
assertVisualResult(result);
```



## OUTPUT

2 of 3 checks passed.

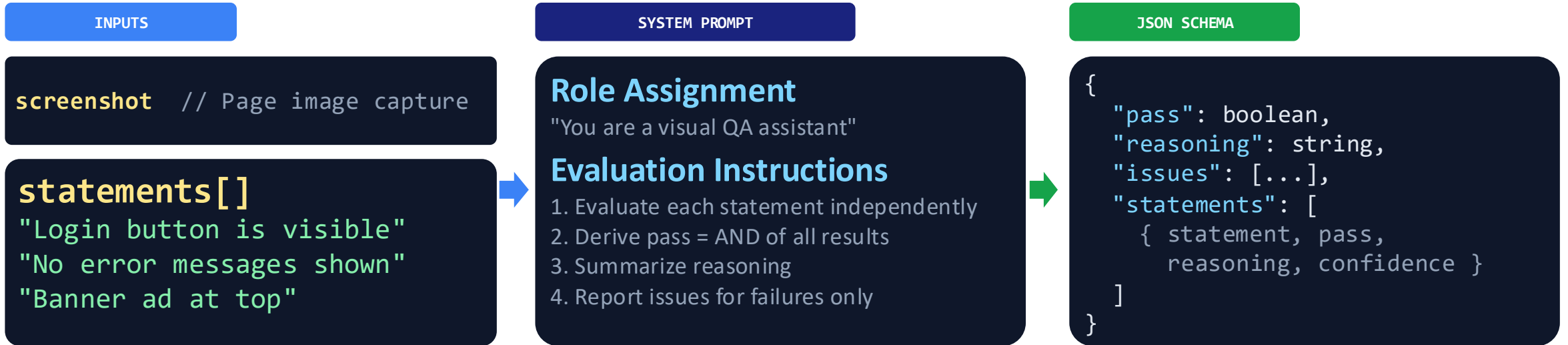
Statements:

- PASS** "The login button is visible" – A green 'Login' button is clearly visible. (high)
- PASS** "No error messages are displayed" – No visible error alerts on the page. (high)
- FAIL** "A banner ad is present at the top of the page" – No banner ad visible. (high)



# Behind the Scenes: System Prompt

How ai.check transforms statements into structured visual QA results



## Key Prompt Concepts

### Role Definition

Sets the LLM as a visual QA assistant that evaluates screenshots precisely and objectively

### Evaluation Instructions

Enforces statements-first evaluation, then derives pass as logical AND of all results

### Confidence Scoring

Each statement gets high / medium / low confidence based on visual clarity and ambiguity

### Structured Issues

Failed statements generate issues with priority, category, description, and suggestion

<https://github.com/nullp2ike/visual-reasoning>

# Integrating Visual Reasoning Checks

diff.e2e.ts

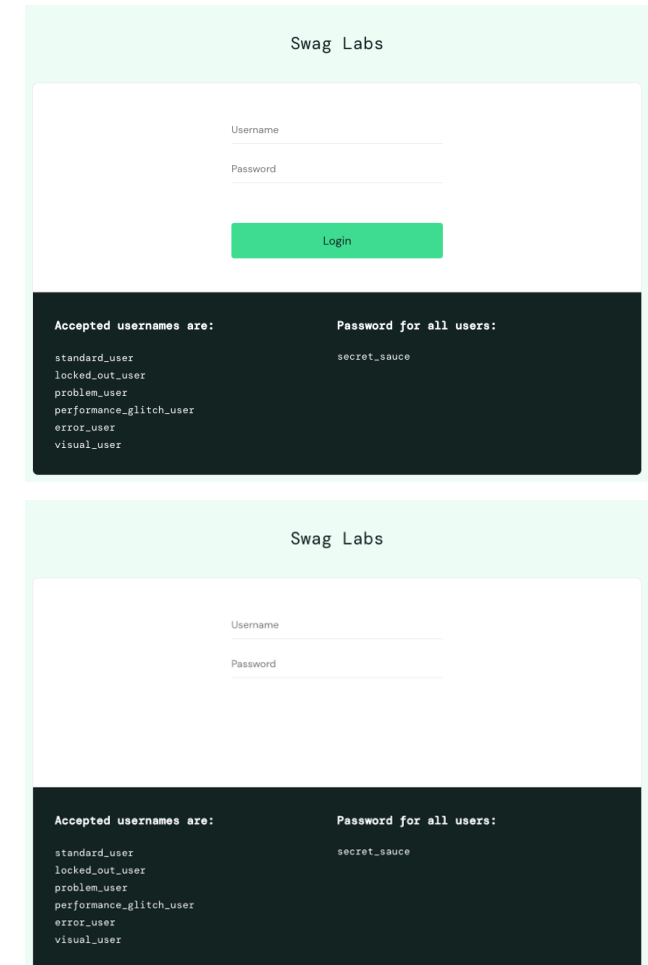
```
const result = await ai.compare(  
  image('before.png'), //with login button  
  image('after.png') //without the login button  
);  
  
expect(result.pass).toBe(true);
```

JSON OUTPUT

```
{  
  "pass": false,  
  "reasoning": "The primary call to action  
  (the Login button) is completely  
  missing in the current version.",  
  "changes": [{  
    "description": "The green 'Login'  
    button has been removed.",  
    "severity": "critical"  
  }]  
}
```

## Visual Regression

Compare baseline vs current with AI



<https://github.com/nullp2ike/visual-reasoning>

# Integrating Visual Reasoning Checks

post-deploy.e2e.ts

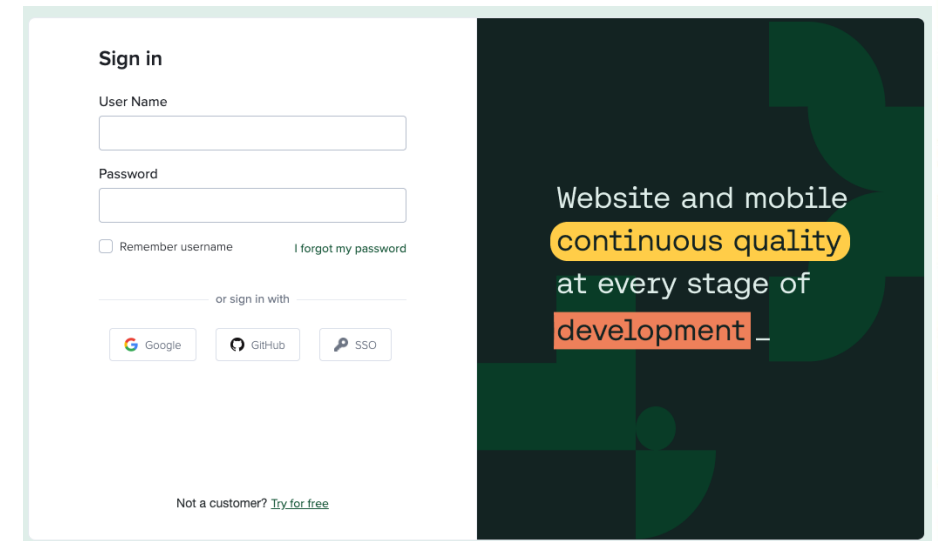
```
const result = await ai.check(screenshot, [  
  'The page is displayed without any critical issues'  
]);  
expect(result.pass).toBe(true);
```

JSON OUTPUT

```
{  
  "pass": false,  
  "reasoning": "The login screen is missing a primary call-to-  
action button (e.g., 'Sign in' or 'Login')...",  
  "issues": [{  
    "priority": "critical",  
    "category": "missing-element",  
    "description": "The primary 'Sign in' or 'Login' button is not  
visible on the form.",  
    "suggestion": "Add a prominent submit button..."  
  }]  
}
```

## In Continuous Delivery Pipelines

Run final visual validation as a post-deploy step



<https://github.com/nullp2ike/visual-reasoning>

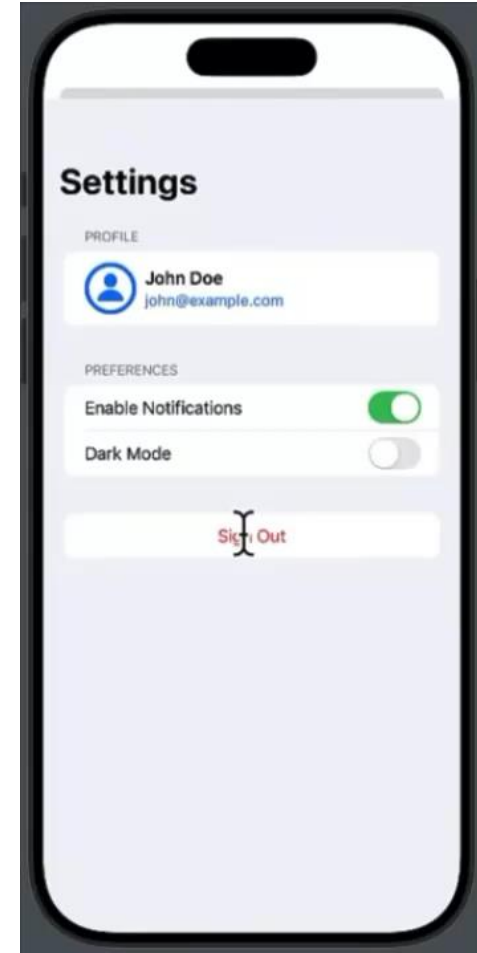
# Video Assertion Support

post-deploy.e2e.ts

```
const result = await ai.check(video, [  
  'You are signed out.'  
]);  
expect(result.pass).toBe(true);
```

JSON OUTPUT

```
{  
  "pass": true,  
  "reasoning": "The requested toast is visible during the sign out flow.",  
  "issues": [],  
  "statements": [{  
    "statement": "The text \"You are signed out.\" is shown in a toast during  
sign out.",  
    "pass": true,  
    "reasoning": "A blue toast displays \"You are signed out.\" at 5.50s,  
visible through 6.50s.",  
    "confidence": "high",  
    "timestampSeconds": 5.5  
  }]  
}
```



<https://github.com/nullp2ike/visual-reasoning>

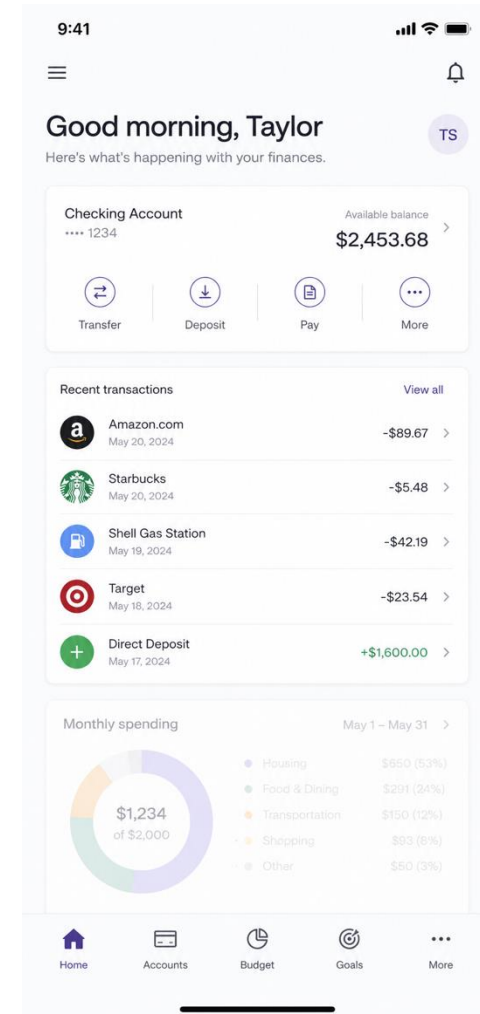
# Accessibility Support

post-deploy.e2e.ts

```
const result = await ai.accessibility(screenshot, {
  checks: [Accessibility.CONTRAST, Accessibility.COLOR_ALONE],
});
expect(result.pass).toBe(true);
```

JSON OUTPUT

```
{
  "pass": false,
  "reasoning": "Most primary content is readable, but some small secondary text and labels appear very light and hard to read.",
  "issues": [{
    "priority": "major",
    "description": "Secondary text and chart labels are low contrast, especially in the monthly spending card.",
    "suggestion": "Darken light gray and pale purple text against the white background."
  }]
}
```



# Privacy & Assertion Library Project Status

## Privacy

### No data collection

The library does not transmit any user data except for what the core functionality does by sending API requests to LLM providers.

**Always verify this yourself before adopting it.**

## Project status

- Proven to be useful in a real project, though **not an official stable release yet**
- Fork it, experiment, and see if it fits your use case.
- Supports currently **OpenAI, Anthropic, and Google** models, but more vendors can be easily added

More on GitHub: [github.com/nullp2ike/visual-reasoning](https://github.com/nullp2ike/visual-reasoning)

# Select Only The Element You Are Interested In

Send the model only the element you care about — sharper answers, lower cost, faster runs.

- **Accuracy** — less visual noise, fewer distractors.
- **Cost** — far fewer image tokens per call.
- **Speed** — less data to upload and process.

```
crop-element.ts
```

```
// WebDriverIO
const cta = await $('~signupCta');
await cta.saveScreenshot('./cta.png');
// Playwright
const cta = page.getByRole('button', { name: 'Sign up' });
await cta.screenshot({ path: 'cta.png' });
// Hand the element image to the model
await ai.check('cta.png', ['Button says Get started']);
```



# Monitoring Production in the Wild

---

Run visual LLM assertions against live sites and apps — catch regressions users would actually see.

- **Synthetic checks** — scheduled crawls hit key flows, screenshot the elements that matter.
- **LLM assertions, not pixel diffs** — ask in plain English: "Is the CTA visible and legible?"
- **Resilient to noise** — shrugs off A/B variants, dynamic content, minor restyles.
- **Catches what users notice** — overlapping text, broken images, cut-off buttons, contrast fails.
- **Page on real failures** — route findings to Slack/PagerDuty with the offending screenshot attached.

monitor-prod.ts

```
// Run every 5 min in prod
await page.goto('https://acme.com');
const hero = page.locator('[data-test=hero]');
await
hero.screenshot({path: 'hero.png'});

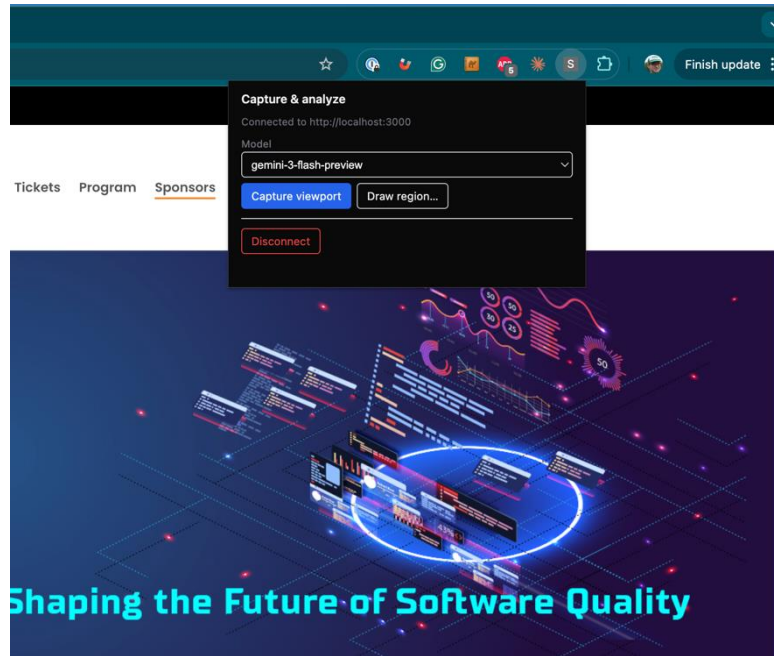
// Visual LLM assertions
const result = await
ai.check('hero.png', [
  'Headline is readable',
  'CTA button visible',
  'No overlap or cutoff',
]);

if (!result.ok) alert(result);
```

# Monitoring Production in the Wild: No-Code Tooling

The screenshot shows a web application interface for 'Sight Checker'. On the left is a sidebar with navigation options: OVERVIEW (Home, Quick Checks), PROJECTS (Test Project 1, https://iststc.com/, Page load checks, Sponsors page, Program page, Tickets page, Landing page), SETTINGS (Global settings, Users, Activity log), and APPEARANCE (Light, Dark, Auto). The main area displays a test configuration for 'Sponsors page' on 'https://iststc.com/'. It includes a 'Steps' list with three actions: 'goto https://iststc.com/', 'click role=button[name="Navigation Menu"]', and 'click body > div:nth-of-type(2) > div > div:nth-of-type(1) > div > ul:nth-of-type(1) > li:nth-of-type(4) > a'. To the right are configuration panels for 'Final capture' (Mode: viewport), 'Auth' (set to None), and 'Tags' (no tags defined). Below is an 'Analysis categories' section with checkboxes for 'visual', 'content', 'functional', and 'accessibility'. A user profile for 'risko.ruus@gmail.com' is visible in the bottom left, and the browser status bar at the bottom indicates 'AdGuard Browser Assistant 1.4.8'.

# Empowering Manual Testing – Browser Extension



**Sight Checker**

OVERVIEW

- Home
- Quick Checks

PROJECTS

- Test Project 1 3 suites
- https://iststc.com/ 1 suite

SETTINGS

- Global settings
- Users
- Activity log

APPEARANCE

Light Dark Auto

---

risko.ruus@gmail.com  
Phase 1 · development

Sign out

Quick Checks

### Quick Check

https://iststc.com/sponsors/

complete viewport gemini/gemini-3-flash-preview \$0.001600 Delete

Capture

Findings

Minor Visual 379x41 at (473,100)

The navigation items and buttons in the header are not vertically centered relative to each other.

nav-alignment-issue

# Empowering Manual Testing – Mobile Companion

Visual checks on real devices — a shortcut, your voice, and an LLM verdict on your phone.

- **Device shortcut** — iOS Shortcut or Android quick-tile fires a screenshot of whatever the tester is looking at.
- **A shortcut for the shortcut** — Bind the shortcut to the phone's action button or to back double/triple tap accessibility features
- **Voice the assertion** — tester speaks the check: "E.g. check the content copy text, critical issues only"
- **Sent for analysis** — screenshot + transcribed prompt are POSTed to the LLM service; verdict returns in seconds.
- **Companion app log** — pass/fail, reasoning, and the captured image stream into a session log, ready to attach to the bug and send to an issue tracker



# Key Takeaways

- 1 LLMs bring semantic understanding to visual testing — they see meaning, not just pixels

# Key Takeaways

- 1 LLMs bring semantic understanding to visual testing — they capture meaning, not just pixels
- 2 Multiple AI providers offer different strengths — use tiered strategies for balancing speed, cost and accuracy

# Key Takeaways

- 1 LLMs bring semantic understanding to visual testing — they capture meaning, not just pixels
- 2 Multiple AI providers offer different strengths — use tiered strategies for balancing speed, cost and accuracy
- 3 LLMs are **non-deterministic**, keep using your existing deterministic tests together with visual AI checks

# Key Takeaways

- 1 LLMs bring semantic understanding to visual testing — they capture meaning, not just pixels
- 2 Multiple AI providers offer different strengths — use tiered strategies for balancing speed, cost and accuracy
- 3 LLMs are **non-deterministic**, keep using your existing deterministic tests together with visual AI checks
- 4 Integration is straightforward — add AI checks as assertions in existing test suites

# Key Takeaways

- 1 LLMs bring semantic understanding to visual testing — they capture meaning, not just pixels
- 2 Multiple AI providers offer different strengths — use tiered strategies for balancing speed, cost and accuracy
- 3 LLMs are **non-deterministic**, keep using your existing deterministic tests together with visual AI checks
- 4 Integration is straightforward — add AI checks as assertions in existing test suites
- 5 Start small, measure impact, and scale gradually across your test infrastructure

**ISTANBUL**  
**SOFTWARE**  
**TESTING**  
**CONFERENCE**

**Thank you!**

**Questions? Get in touch:**

[risko.ruus@gmail.com](mailto:risko.ruus@gmail.com) | [linkedin.com/in/risko-ruus/](https://www.linkedin.com/in/risko-ruus/)