

Lessons Learned from AI-Powered Visual Reasoning Feedback

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

Risko Ruus

June 5, 2026



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

Traditional UI Assertions



Agenda

Traditional UI Assertions

**Visual AI Reasoning: Core
Concepts & Basic Usage**



Agenda

Traditional UI Assertions

**Visual AI Reasoning: Core
Concepts & Basic Usage**

**Advanced Analysis
Techniques**



Agenda

Traditional UI Assertions

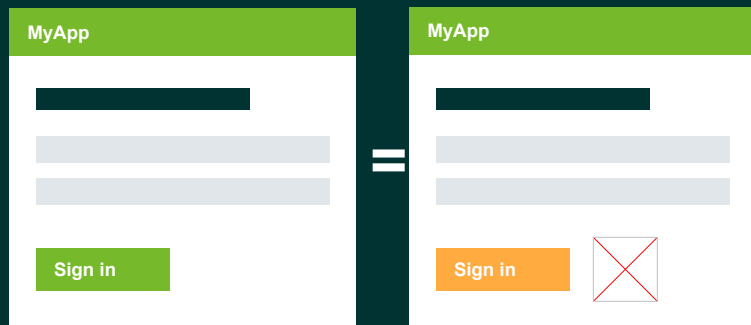
**Visual AI Reasoning: Core
Concepts & Basic Usage**

**Advanced Analysis
Techniques**

**Integrating Visual
Reasoning Checks into
Workflows**



Traditional UI Assertions – Full Match



Version 1

Version 2

Problems:

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

High Maintenance

Brittle



Example Diff Output

Match level simulation
Strict

Set region

Baseline
Strict | iOS 16.1 | Safari 16.1 | 375x812 | iPhone X

Checkpoint
Strict | iOS 16.1 | Safari 16.1 | 375x812 | iPhone X

FRNTR

Names, SKUs, categories

Shop

Let Nature Into Your Living Room

Choose from a wide variety of plants that add life and style to any space.

Japandi Interior Design

Japandi interior design is a hybrid of Japanese and Scandinavian styles. The style is increasingly popular here to stay. The style creates in

Introducing Our New Sale

Choose from a wide variety of plants that add life and style to any space.

Japandi Interior Design

Japandi interior design is a hybrid of Japanese and Scandinavian styles. The style is increasingly popular here to stay. The style creates in

Alocasia Portodora

€232.00

€208.80 -10%

Monstera Deliciosa

€24.00

€20.00 -17%

Alocasia Portodora

€8.80 -10%

Your price

Monstera Deliciosa

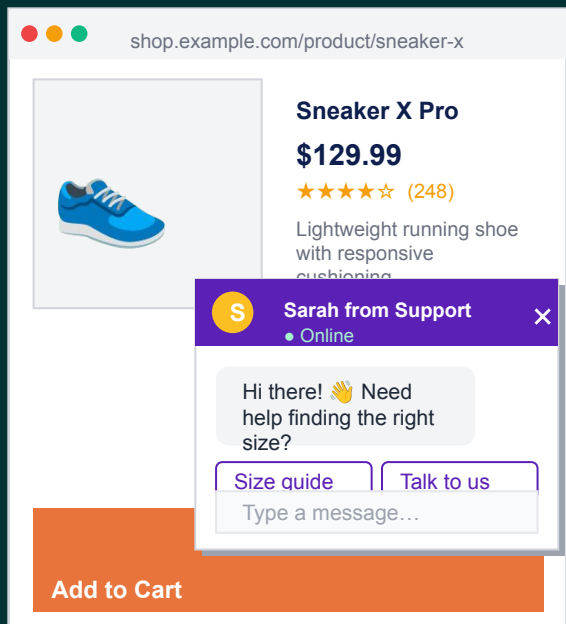
€24.00

€20.00 -17%



Real-World Examples: When isDisplayed/isVisible is Not Enough

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



shop.example.com/product/sneaker-x

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

Sarah from Support
● Online

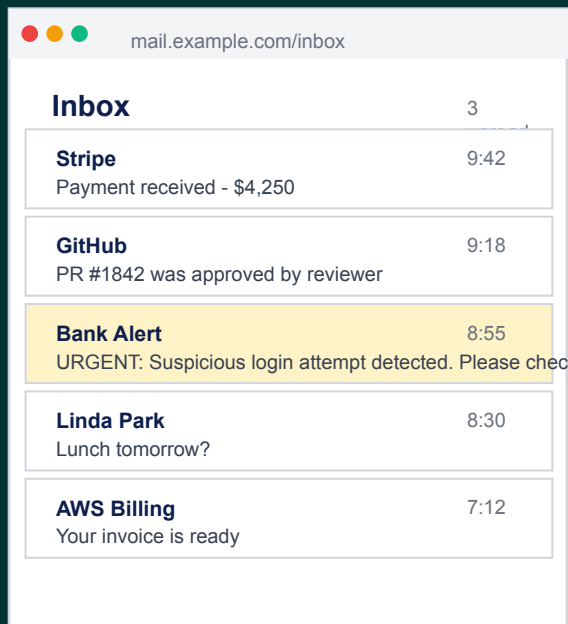
Hi there! 🙋 Need help finding the right size?

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

Type a message...

[Add to Cart](#)

Covered by overlay



mail.example.com/inbox

Inbox 3

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

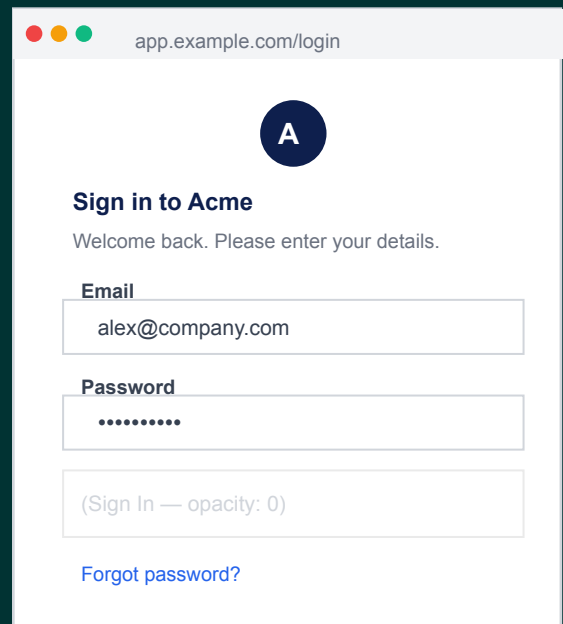
GitHub 9:18
PR #1842 was approved by reviewer

Bank Alert 8:55
URGENT: Suspicious login attempt detected. Please check your account details.

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)

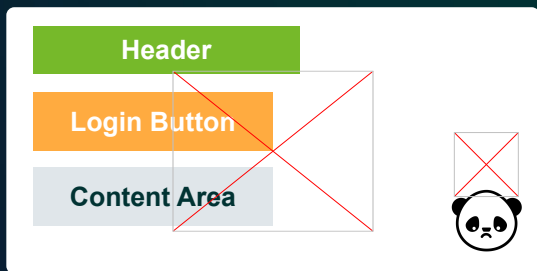
[Forgot password?](#)

Opacity: 0

Traditional UI Assertions – Specific Element Check

Aka - The Giant Panda Problem

Goal - Verify that a login button is displayed



PASS



MISSED!

Problems:

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

Blind Spots

False Confidence



Visual AI Reasoning: Core Concepts & Basic Usage

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 reasons

The LLM breaks the image and prompt text into tokens. Then it spends any internal thinking tokens and finally generates a response.

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*



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



Sauce Labs Backpack

carryallTheThings() 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, 1AAA 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

Analyze this image and identify all UI, content, and functional issues.

Severity levels:

Critical – Blocks the main flow/actions, damages trust

Major – Clear defect, workaround exists

Minor – Cosmetic, unlikely to affect users

Output valid JSON using this schema

→

SYSTEM PROMPT



```
{
  "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..."
  }
}
```

JSON SCHEMA



Example Vision LLM Prompt Output

JSON OUTPUT

```
{
  "id": 1,
  "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

Swag Labs



1

Sauce Labs Backpack

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

\$18.24

Add to cart



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

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 sleeves and bottom won't unravel.

\$14.83

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 Fleece Jacket

It's necessary 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



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 Add to cart

6

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

Detects	Pixel-level differences
Reports	Changed vs unchanged pixels
False Positives	High
Feedback	"Pixels differ at (x,y)"
Exclusions	Manual exclusion areas per image
Speed	Milliseconds (local compute)
Cost	Near zero (local CPU)
Determinism	Deterministic

AI-Powered Visual Reasoning

Semantic UI issues
Layout, hierarchy, meaning, accessibility
Medium (contextual reasoning)
"Button is hidden by overlay"
Ignore prompts
Seconds per image
API token fees per invocation
Non-deterministic!

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



Advanced Analysis Techniques – Model Selection Overview

OpenAI	Anthropic	Google	Open Weights
<ul style="list-style-type: none">• GPT-5.5• GPT-5.4 Pro• GPT-5.4• GPT-5.4 Mini• GPT-5.4 Nano	<ul style="list-style-type: none">• Opus 4.7• Sonnet 4.6• Haiku 4.5	<ul style="list-style-type: none">• Gemini 3.5 Flash• Gemini 3.1 Pro• Gemini 3.1 Flash-Lite• Gemini 3 Flash	<ul style="list-style-type: none">• Qwen 3.5 397B A17B• Qwen 3.5 27B• Kimi K2.6• ...

Thinking effort is configurable on most models

A **"reasoning budget"** setting (minimal / low / medium / high, or a token cap) controls how many internal reasoning tokens the model spends before answering.

Higher effort = better accuracy on hard tasks at more latency & cost; lower effort = faster, cheaper responses.

**as of: May 20, 2026*



Visual Reasoning Intelligence (MMMU Pro evaluation)

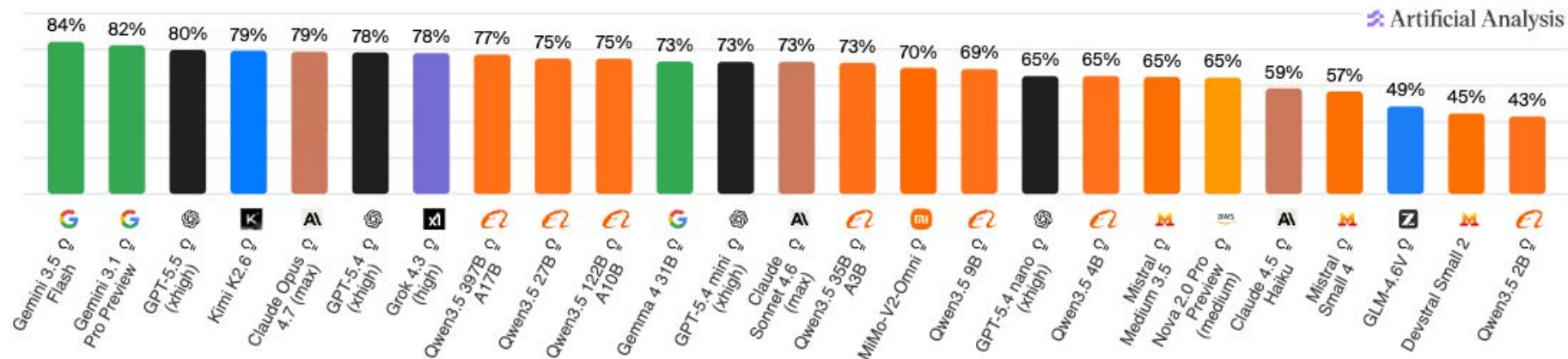
25 of 190 models



0 of 417 model & provider com...



Visual reasoning intelligence: MMMU Pro evaluation



💡 Reasoning models are indicated by a lightbulb icon.

MMMU Pro

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

*as of: May 20, 2026

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

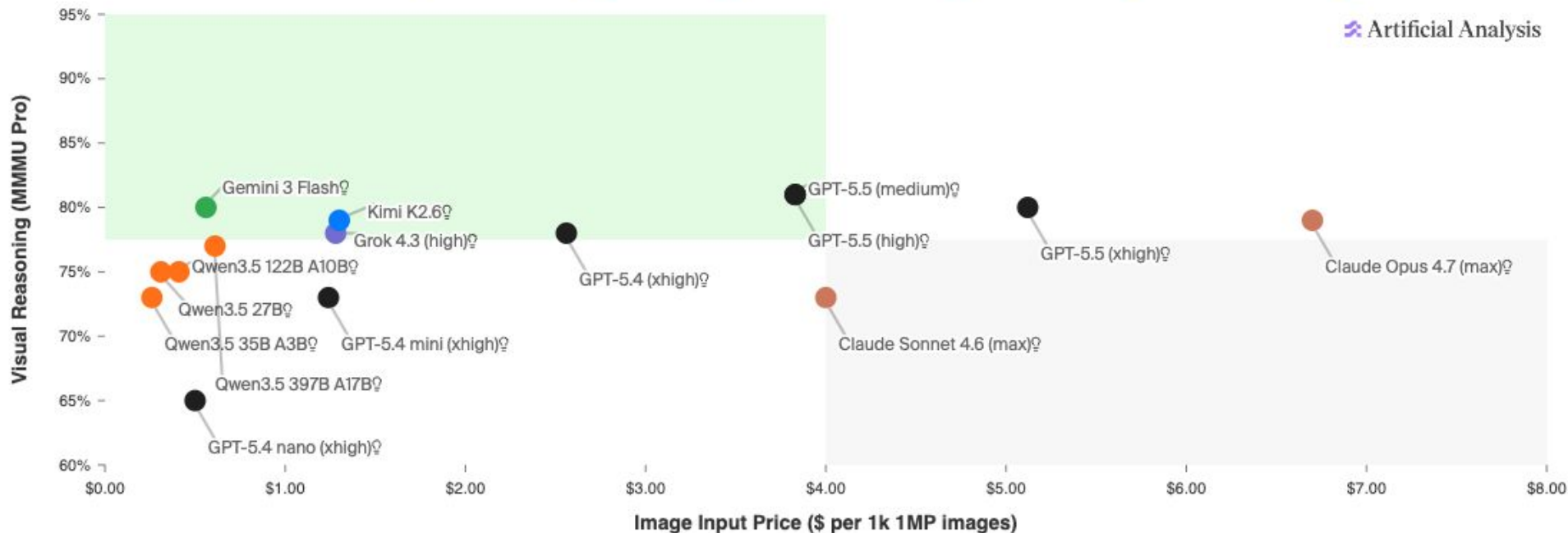
Visual Models Comparison – Visual Reasoning vs. Price*

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

Artificial Analysis



*as of: May 20, 2026

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

Visual 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>

Dealing with Duplicate Findings

The problem

- AI tools flag the same bug with different wording each run
- Reviewers waste time closing duplicates instead of fixing bugs
- Simple text matching misses rephrased findings

We need to compare by meaning, not just words. An **embedding** model can be used for that

Example: same bug, different words

Run 1

“Login button on mobile is cut off; users can’t tap submit.”

Run 2

“Sign-in CTA overflows viewport on iPhone — primary action not clickable.”

The second run finding gets linked to the initial one — no duplicate ticket created.

✓ **90% match — treated as a duplicate**



Agent Review: Cross-Checking VLM Findings

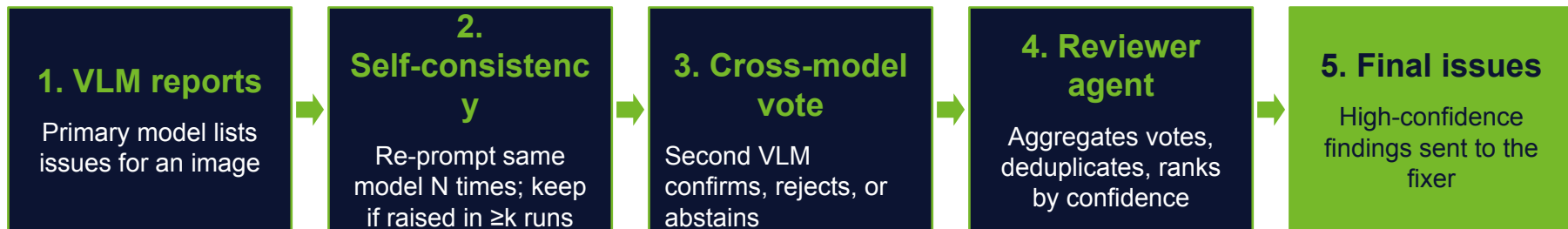
The problem

- An LLM may hallucinate or flag subjective issues
- Hard to separate signal from noise in one pass
- Acting on every report wastes engineering time

Agent-based review

- Reviewer agent rechecks each reported issue
- Self-consistency: N calls to the same model
- Cross-model: a second VLM votes
- Keep issues confirmed above a threshold

How the review loop works



Integrating Visual Reasoning Checks

```
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"
]);
assertVisualResult(result);
```

OUTPUT

2 of 3 checks passed.

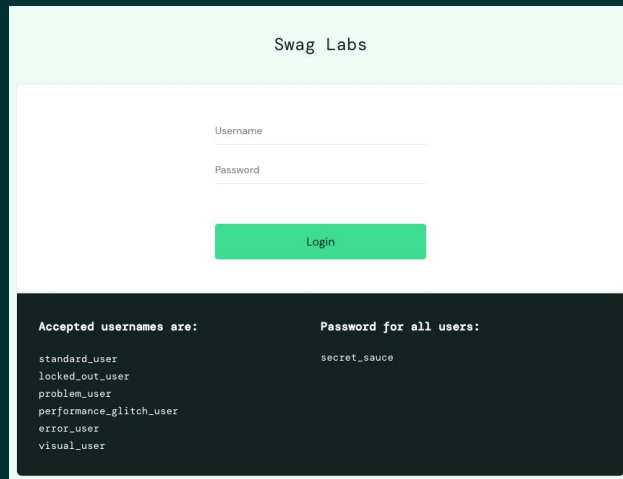
PASS "The login button is visible" – green Login button clearly visible (high)

PASS "No error messages are displayed" – no visible error alerts (high)

FAIL "A banner ad is present at the top" – no banner ad visible (high)

Drop-in for any framework

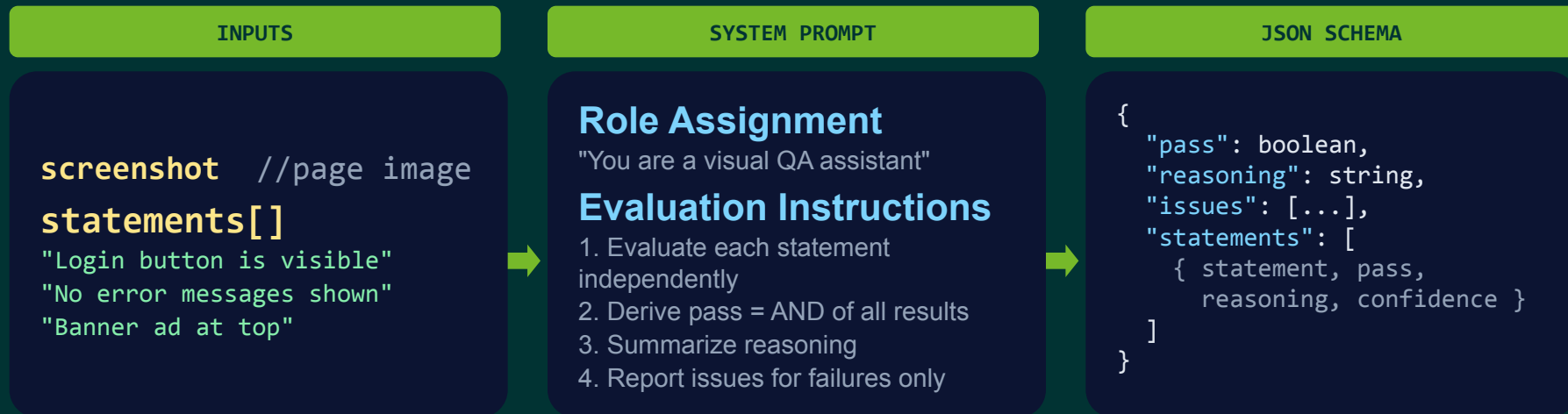
Selenium • Playwright • Cypress • Appium



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

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 evaluating screenshots precisely

Evaluation Instructions

Statements first, then derive pass as logical AND of all results

Confidence Scoring

High / medium / low confidence per statement based on visual clarity

Structured Issues

Failed statements emit priority, category, description, suggestion

Visual Regression with AI

diff.e2e.ts

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

JSON OUTPUT

```
{  
  "pass": false,  
  "reasoning": "Login button is missing in the  
current version",  
  "changes": [{  
    "description": "Green Login button removed",  
    "severity": "critical"  
  }]  
}
```

BEFORE

Swag Labs

Username

Password

Login

Accepted usernames are:

standard_user
locked_out_user
problem_user
performance_glitch_user
error_user

Password for all users:

secret_sauce

AFTER

Swag Labs

Username

Password

Accepted usernames are:

standard_user
locked_out_user
problem_user
performance_glitch_user
error_user
visual_user

Password for all users:

secret_sauce

NTD



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



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

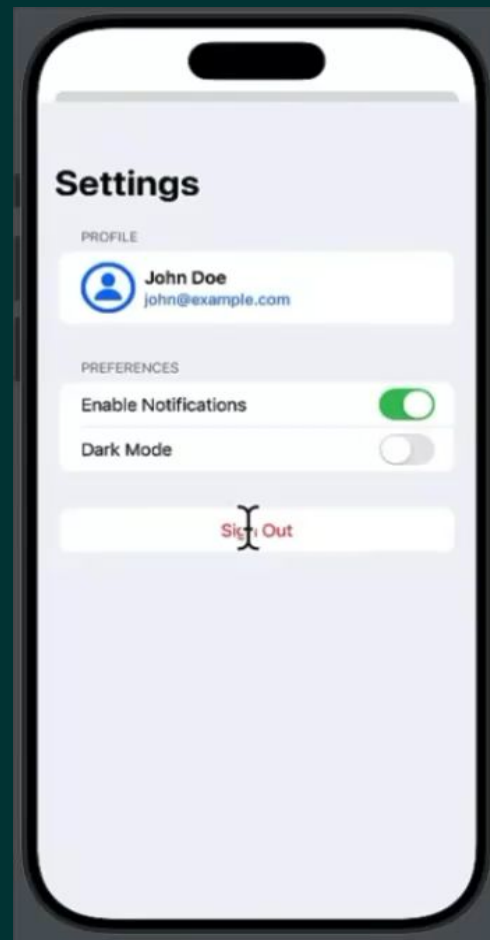
Video Assertion Support

post-deploy.e2e.ts

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

JSON OUTPUT

```
{  
  "pass": true,  
  "reasoning": "Toast is visible during sign out flow.",  
  "issues": [],  
  "statements": [{  
    "statement": "\"You are signed out.\" shown in toast",  
    "pass": true,  
    "confidence": "high",  
    "timestampSeconds": 5.5  
  }]  
}
```



Accessibility Support

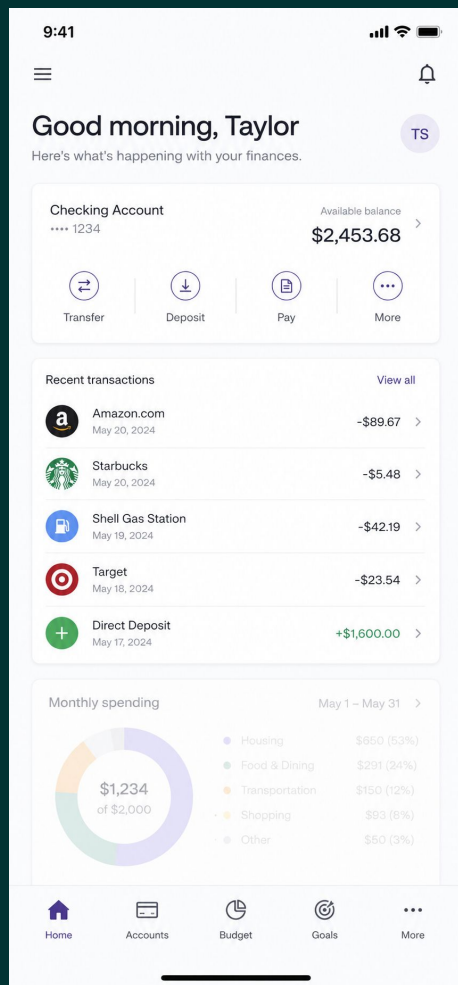
a11y.e2e.ts

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

```
expect(result.pass).toBe(true);
```

```
{
  "pass": false,
  "reasoning": "Primary content readable; some
  small labels appear too light to read.",
  "issues": [{
    "priority": "major",
    "description": "Secondary text and chart labels
    low contrast on the spending card.",
    "suggestion": "Darken light gray and pale purple
    text against the white background."
  }]
}
```

JSON OUTPUT



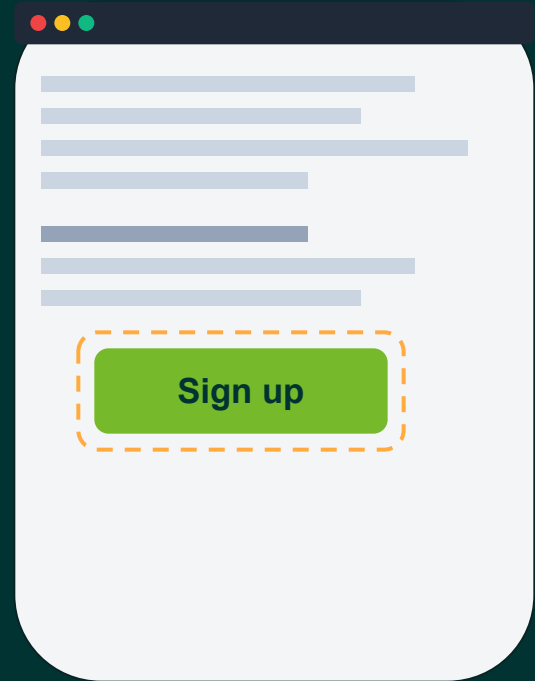
Select Only the Element You Care About

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']);
```



Crop → only this element reaches the model



Sight Checker

OVERVIEW

- Home
- Quick Checks

PROJECTS

- NTD 2026 1 suite
 - Partner's page 1 test
 - Navigation from homepa...

SETTINGS

- Global settings
- Users
- Activity log

APPEARANCE

Light Dark Auto

NTD 2026 / Partner's page

Navigation from homepage to partner's page

mobile 390x844 https://nordictestingdays.eu

Rename Delete Clone Record Run

Editor History Open issues (2) FP suppressions

Steps

4 steps, ending with a viewport capture.

Extract action group... Add step

1	goto	https://nordictestingdays.eu	ENTRY	
2	click	role=button[name="Toggle navigation"]		Edit Delete
3	click	role=link[name="Partners"]		Edit Delete
4	capture	viewport		Edit Delete

Final capture Edit

Mode: viewport

Auth

— None —

Tags

No tags in this project's vocabulary yet — create some on the project Tags settings page.

Apply to test

Analysis categories

Override the project's active category set for this test only.

Override project categories Inheriting: visual, content, functional

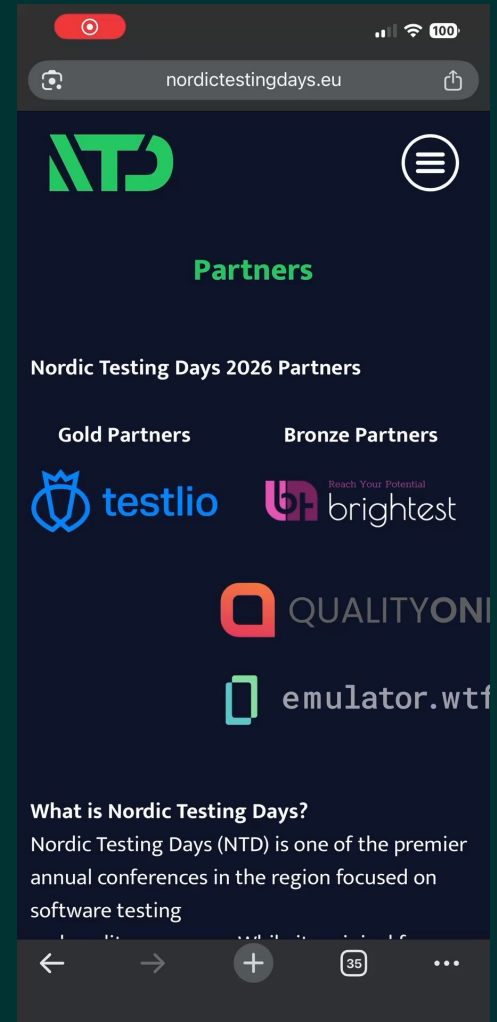
visual content functional accessibility

Save

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.
- **Shortcut for the shortcut** — bind it to the phone's action button or to a back double/triple-tap gesture.
- **Voice the assertion** — tester speaks the check: "check the content copy, critical issues only."
- **Sent for analysis** — screenshot + transcribed prompt 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 be submitted as a bug.



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 to balance speed, cost, and accuracy

3

LLMs are **non-deterministic** — keep deterministic tests alongside visual AI checks

4

Integration is straightforward — add AI checks as assertions in existing test suites

5

Start small, measure impact, and don't be afraid to build tools



Thank you

Questions? Let's keep the conversation going.

GET IN TOUCH



EMAIL

risko.ruus@gmail.com



LINKEDIN

[linkedin.com/in/risko-ruus](https://www.linkedin.com/in/risko-ruus)

