

The AI-first WordPress site

CRAWLER TO CITATION

Alain Schlessler



// WCEU / 2026
// KRAKOW · JUN 06

SATURDAY · 2:45 PM CEST
DEVELOPMENT · TRACK 1

CRAWLER

TO CITATION.



> ABSTRACT

The AI-first WordPress site. Make it readable, citable, and measurable by the machine layer that just appeared on top of the web.

> AUTHOR

ALAIN SCHLESSER Agentic architect & engineer · WP-CLI maintainer · Google Developer Expert
alainschlessers.com

// SECTION DIVIDER
// ACT ONE OF THREE

5 MIN
5 SLIDES

ACT I.

THE PROBLEM


▶ PROBLEM
▶ 5 MIN

AI is now a traffic source. **It sends visitors, and it takes content.**
The shift is settled — we're here to talk about what to do.

// ACT I · THE PROBLEM

// THE OPENER

PICTURE THIS



Sometime in the last week, an AI system decided whether to recommend your site as the answer to someone's question.

**YOU WEREN'T
IN THE ROOM.**

AND MOST OF THE TIME, THE ANSWER WAS NO.



— A.S.

// ACT I · THE PROBLEM

// FIGURE 03.1

THE AUDIENCE JUST DOUBLED.

THE SHIFT, IN ONE NUMBER

— PAGE 14 · FOLD C —

**BOTS
OUTNUMBER
HUMANS**

▶ READ THIS

▶ AI BOT TRAFFIC · YOY · OPEN WEB

+187%

Growth across the open web, Jan - Dec 2025. The direction is one-way. Don't argue with it - design for it.

— CLOUDFLARE RADAR · Q1 2026

Humans, same window **+3.1%**

// ACT I · THE PROBLEM

// FIGURE 04.1

QUALITY, NOT JUST VOLUME

A FULL

REVERSAL.

BEFORE → AFTER

in 12 months.

► MARCH 2025

AI traffic converted -38 % vs non-AI.
The skeptic position made sense.

-38%

► MARCH 2026

AI traffic converts +42 %, RPV +37 %,
time on site +48 %, pages/visit +13 %

+42%



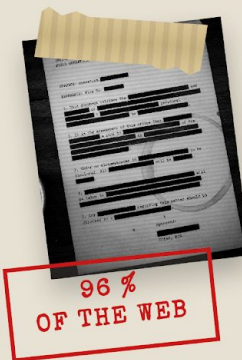
// ACT I · THE PROBLEM

// THE GAP

WHERE MOST WP SITES ARE

THREE WAYS

TO BE INVISIBLE.



INVISIBLE

No robots.txt strategy. No schema. No measurement. The AI ecosystem flows past you in both directions.

1

VULNERABLE

Default-open to training crawlers. No governance.

2

Content harvested by anyone. No idea who, or how much.

OLD-PARADIGM

Optimised hard for Google 2018. Nothing ready for the AI layer that just appeared on top of it.

3

// SECTION DIVIDER
// ACT TWO OF THREE

18 MIN
16 SLIDES · 5 LAYERS

ACT II.

THE STACK

▶ STACK
▶ 5 LAYERS

Five layers — **crawlers · robots · schema · content · measurement.**

Take notes on the ones you don't have yet.



// ACT II · STACK · LAYER 1

// CRAWLER AWARENESS

THREE BOTS.

THREE STANCES.

THE THREE TIERS

- tip:

agents = users.
treat them so.

> TIER 01

TRAINING

GPTBot · ClaudeBot · CCBot
Google-Extended · Bytespider
· Amazonbot

BLOCK OR RATE-LIMIT
UNLESS PAID

> TIER 02

SEARCH / CITATION

OAI-SearchBot ·
PerplexityBot
ChatGPT-User · AppleBot-
Extended · Bingbot

ALLOW.
THESE DRIVE CITATIONS.

> TIER 03

AGENT

OpenAI Operator
Anthropic Computer Use
Perplexity Comet

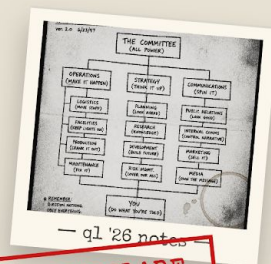
ALLOW.
TREAT LIKE A USER.

// ACT II · STACK · LAYER 1
// VENDOR SHARE, Q1 2026

MULTI-VENDOR REALITY

89 %

IS TRAINING.



— q1 '26 notes —
CLOUDFLARE
Q1 2026

1 **GOOGLEBOT**
31.6% — the elephant. Still leads, no longer a monopoly.

3 **GPTBOT**
12.0% — OpenAI training. Clear opt-out.

5 **APPLEBOT-EXTENDED**
5.8% — up 124% in Q1. Siri Intelligence is here.

2 **META-EXTERNALAGENT**
16.7% — quietly enormous. Mixed-purpose crawler.

4 **CLAUDEBOT**
11.7% — Anthropic. Distinct UAs for training vs search.

6 **THE REST**
22.2% — Bytespider, Amazonbot, CCBot, PerplexityBot, dozens more.

// ACT II · STACK · LAYER 1

// WHO'S AT YOUR DOOR

YOU CAN SEE

ALL OF THEM.

TOOLS YOU ALREADY HAVE

— FIELD GUIDE

**WHO'S
AT YOUR DOOR**

USER-AGENT

Every major vendor publishes theirs. UA-sniffing works — until it doesn't. Trivially spoofable.

1

IP / ASN

OpenAI, Anthropic, Apple publish ranges. UA + ASN = verified bot.

2

SERVER LOGS

One grep one-liner on your access log. Visibility today, no plugins.

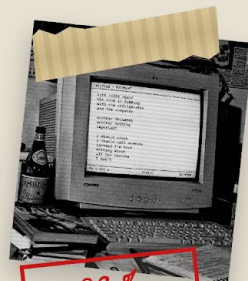
3

// ACT II · STACK · LAYER 2
// ROBOTS.TXT STRATEGY

DEFAULT

IS NOT A STRATEGY.

WHAT WP SHIPS



96 %
OF SITES

ROBOTS.TXT

WORDPRESS 6.5 DEFAULT

```
# wordpress default - same since 2008
```

```
User-agent: *
```

```
Disallow: /wp-admin/
```

```
Allow: /wp-admin/admin-ajax.php
```

```
# zero ai-specific directives.
```

```
# every ai crawler gets exactly the same treatment.
```

```
# this is not a stance - it's the absence of one.
```

```
// ACT II · STACK · LAYER 2  
// THE STRATEGIC CONFIG
```

THREE TIERS.

THREE DECISIONS.



ROBOTS.TXT

THE 2026 CONFIG

```
# tier 1 - training (default: deny)  
User-agent: GPTBot  
User-agent: ClaudeBot  
User-agent: CCBot  
Disallow: /  
  
# tier 2 - citation (default: allow)  
User-agent: OAI-SearchBot  
User-agent: PerplexityBot  
Allow: /  
  
# tier 3 - agent (allow, like users)  
User-agent: ChatGPT-User  
Allow: /
```

YOUR CONTENT.

YOUR DECISION.



THE RETURN IS ZERO

1

Training crawlers take, don't give. No citation, no referral, no attribution. Bandwidth on the way in is your bill.

LEGAL DIRECTION

2

Publishers vs major AI labs in active litigation. Direction of travel: more control, not less. A clear opt-out today is cheap insurance.

TRACK RECORD

3

Some operators have a documented history of ignoring opt-outs. The compute they spend crawling you is bandwidth you pay for.

RE-AUDIT QUARTERLY

4

The line between training and search is shifting. Some vendors run unified crawlers. Not a permanent answer — a current one.

// ACT II · STACK · LAYER 2
// WHEN HONOR-SYSTEM FAILS

WHEN HONOR

SYSTEM FAILS.

ENFORCEMENT LAYERS



belt &
braces.
both.

WAF RULE

Cloudflare AI Crawl Control, Fastly's equivalent, Akamai's. One toggle in a dashboard.

1

RATE-LIMIT + IP BLOCK

Server-level. More granular. Useful when you want to allow a vendor sometimes, not always.

2

WEB BOT AUTH

HTTP message signatures. Cryptographic identity for well-behaved bots. Standards work in progress.

3

// ACT II · STACK · LAYER 3

// THE REFRAME

SCHEMA IS THE DATA LAYER



Schema used to be an SEO trick. In 2026 it's something else.



TREAT IT LIKE AN API. THE REST OF THIS LAYER MAKES SENSE.

— A.S.

THREE

INDEPENDENT RECEIPTS.



PLATFORM

1 Fabrice Canel,
Principal PM Bing,
SMX Munich: "Bing &
Copilot use schema
to feed their LLMs."
Not me inferring – a
Microsoft PM on a
public stage.

TRAINING TIME

2 LLM pre-training
pulls from Common
Crawl. Your
structured data is
in that snapshot.
Every model inherits
the structure.
Effect is one-way
and permanent.

PROTOCOL LAYER

3 NLWeb, MCP-based
approaches, llms.txt
– several proposals
worth taking
seriously. All of
them speak schema.
Disclosure: I led
NLWeb work at Yoast.

STUB SCHEMA HURTS YOU.

DO IT
PROPERLY

► CITATION RATE · MINIMAL / GENERIC SCHEMA

41.6%

Pages with no schema were cited 59.8% of the time. Rich attribute-complete schema hit 61.7%. Half-finished schema looks like low effort. AI systems notice.

— GROWTH MARSHAL · N=730

Either commit, or don't bother. **61.7%**

// ACT II · STACK · LAYER 3

// THE @ID GRAPH

GRAPH BEATS BLOBS

CONNECT

EVERY @ID.

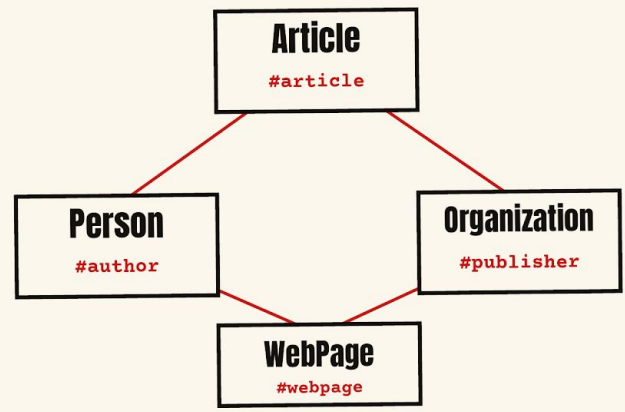
graph > blob.
every. single.
time.

► BEFORE — isolated blobs



→ AI sees three loose objects.

► AFTER — interconnected @id graph



// ACT II · STACK · LAYER 4

// CONTENT PATTERNS

EVIDENCE-BASED

THREE EDITS.

+30 - 40 %.

// content edit pass
cite. quote. count.
→ ~30-40 % lift

PRINCETON
GEO-BENCH

CITE SOURCES

Pages that cite get cited. Add references.

- 1 Hyperlink to primary sources. Position-adjusted lift: ~30-40 %.

QUOTE AUTHORITIES

Pull direct quotes into blockquotes. AI treats your page as quoting an authority. Same lift band.

- 2

USE STATISTICS

Replace "a lot" with the actual number. Specific figures get pulled when AI grounds its own claims.

- 3

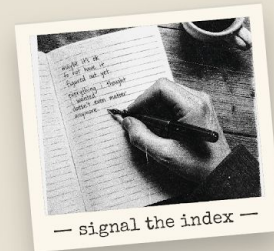
// ACT II · STACK · LAYER 4

// THE TWO HEURISTICS

FRESH BEATS

FOSSIL.

FRESHNESS + ENTITIES



FRESHNESS SIGNAL

Microsoft Canel: "Generative AIs value fresh content — as a reference check of LLM training."

1

His channel:

INDEXNOW

. WP plugins exist. Cheap.

ENTITY DENSITY

Named entities + internal links improve AI comprehension. King (iPullRank): relevance, not authority, drives AI Overview placement.

2

// ACT II · STACK · LAYER 4

// AWKWARD TRUTH

BECOMING THE SUBSTRATE

WIKIPEDIA.

REDDIT. YOU.



SEMRUSH
100M+ CITES

CHATGPT

- 1 Wikipedia 26-48 % of top-10 share.
Reddit went 60 % → 10 % in six weeks
after a Google parameter change.

GOOGLE AI MODE

- 2 Reddit, LinkedIn, YouTube growing.
Medium, Quora declining.

PERPLEXITY

- 3 Reddit at ~40 %. WSJ, NYT, Bloomberg
absent from top 20.

THE PLAY

- 4 You can't out-Reddit Reddit. Be the
original-research source. The thing
the forum thread links to.

REFERRERS

TO WATCH.



CHAT.OPENAI.COM

- 1** The biggest single source for most categories.

PERPLEXITY.AI

- 2** Smaller but converts hard. Power users.

GEMINI.GOOGLE.COM

- 3** Growing fast. Tied into the Google account flow.

COPILOT.MICROSOFT.COM

- 4** Enterprise traffic. Edge default.

X.COM (GROK)

- 5** Stripped sometimes. Show up as x.com.

KAGI.COM

- 6** Small, paid, high-intent.

// ACT II · STACK · LAYER 5

// MINIMAL DASHBOARD

FOUR METRICS.

NO MORE.

FOUR METRICS

start small.

measure something.
beats nothing.

VISITS BY AI REFERRER

- 1 Undercounted baseline. Trend over time is what matters, not absolute.

CONVERSION / ENGAGEMENT BY REFERRER

- 2 Pages, time, goal completions. Does the Adobe pattern hold for you?

CITATION SHARE

- 3 Pick ten queries. Rotate weekly. Log which sites get cited and track your share over time.

CRAWLER SHARE

- 4 Grep your access log weekly. Training-tier vs citation-tier bots. The bandwidth picture you don't see in GA.



Schema is the foundation. Several proposals are competing for what gets built on top — and we don't yet know which one wins.

**THE NEXT LAYER ISN'T
ABOUT BEING READABLE.
IT'S ABOUT BEING **CALLABLE.****



// SECTION DIVIDER
// ACT THREE OF THREE

4 MIN
4 SLIDES

ACT III.

► ACTION
► TONIGHT

DO THE WORK



Two timelines. **Tonight**, and the next thirty days.
You can be ahead of 96 % of the web before you sleep.

DO THIS

TONIGHT.



▶ ACTION
▶ ITEMS

AUDIT YOUR ROBOTS.TXT

- 1** Explicit decision per tier: training, citation, agent.

ATTRIBUTE-COMPLETE SCHEMA

- 3** Organization. Person. Article. Stub schema underperforms none.

EDIT FIVE CORNERSTONE PAGES

- 5** Citations + quotes + statistics. Princeton GEO patterns.

VERIFY YOUR SCHEMA GRAPH

- 2** @id cross-references. Not isolated JSON-LD blobs.

RUN A READINESS SCANNER

- 4** One scan today. Capture the baseline score.

ADD LLMS.TXT + REFERRER TRACKING

- 6** Five minutes each. Low risk. Cheap signal.

FOUR WEEKS.

FOUR THEMES.

— CALENDAR B
**FOUR
WEEKS**

WEEK 01 — INSTRUMENT

- 1** Wire the dashboard. Baseline referrals, citations, crawler share. Measurement is slowest to start, so front-load it.

WEEK 02 — SCHEMA

- 2** Top 10 pages to full attribute-completeness. Don't add types — fill in what's there.

WEEK 03 — CONTENT

- 3** Top 5 cornerstone articles. Citations, quotations, statistics. Five well-done beats fifty half-done.

WEEK 04 — REVIEW & PUBLISH

- 4** Iterate where the signal is. Publish your AI-readiness page on your site. Public commitment for round two.

// ACT III · TAKE IT HOME

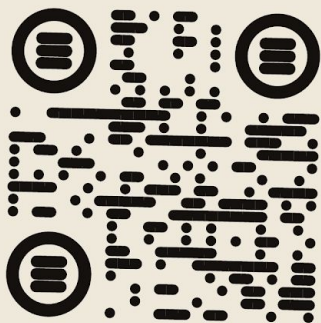
// INTERACTIVE GUIDE

AUDIT YOUR OWN SITE

RUN THE

CHECKLIST.

> SCAN TO START



tiki.tf/wceu2026

> INTERACTIVE GUIDE

Walk every **layer**
on your **own site**,
one tick-box at a time.

Crawlers. Robots. Schema.
Content. Measurement.

– start tonight.

// THE END
// THANK YOU + Q&A

FIN.
OVER TO YOU

THANK YOU. QUESTIONS?

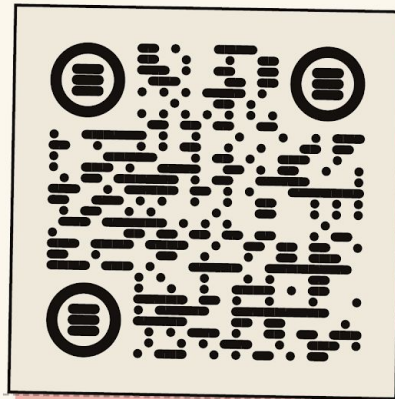
► YOUR TURN

Ask me **anything**.

Pick a **layer**,
share what's **working**,
or push back on a **claim**.

– **let's talk.**

► SCAN TO CONNECT



@schlessera · alainschlessor.com



WordCamp
Europe
Kraków 2026

THANK YOU!

