Technical Analysis • February 2026

# Marking AI Content: A Compliance Guide for EU AI Act Article 50

The EU AI Act (Article 50) defines transparency obligations for AI providers and requires providers of AI-generated content to mark outputs in a machine-readable format that is effective, interoperable, robust, and reliable. With enforcement set to begin in August 2026, this article studies the current landscape and lays out a path to compliance.

| Author | Published | Reading time |
|---|---|---|
| Vibhu Ganesan | February 2026 | 18 minutes |

## Key Takeaways

**1** **Single-layer marking is not sufficient.** The Draft Code is explicit: "no single active marking technique suffices." Providers relying on C2PA metadata alone fail robustness. Those using proprietary watermarks alone fail interoperability.

**2** **No major provider is fully compliant.** Our assessment of ten providers—including Google, OpenAI, Meta, and Microsoft—finds gaps across the board. The industry has not yet implemented the multi-layer marking architecture that the Draft Code of Practice requires.

**3** **The technical path forward is becoming clearer.** C2PA v2.3's soft binding specification enables hybrid architectures that satisfy both robustness and interoperability. The building blocks exist—providers must assemble them.

**4** **Detection access is still overlooked.** Most providers have focused on marking but not on enabling third-party verification at scale. The Draft Code requires "free-of-charge interfaces" that few providers currently offer.

**5** **Privacy requirements are entangled with transparency requirements.** Marking AI-generated content shouldn't violate users' right to privacy. While the Draft Code doesn't explicitly emphasize data minimization in generation of metadata, privacy preservation should be a key design consideration.

**6** **Six months remain.** Providers without watermarking infrastructure face the greatest implementation risk. Watermark integration typically requires six to twelve months—leaving no margin for delay.

# The compliance mandate

Article 50 of the EU AI Act establishes among the first binding requirements for AI content marking in a major jurisdiction. The stakes are significant: non-compliance carries penalties of up to €15 million or 3% of global annual turnover.

The regulation requires that AI-generated content be "marked in a machine-readable format and detectable as artificially generated or manipulated." Four criteria govern what counts as adequate marking: it must be effective, interoperable, robust, and reliable. In December 2025, a Draft Code of Practice was released and this Code of Practice translates these criteria into operational requirements.

The Code's most consequential clarification is that multi-layer marking is mandatory. The reasoning is straightforward. Metadata-only approaches (such as C2PA) fail robustness because metadata is stripped by screenshots, social media uploads, and format conversions—precisely the channels through which AI

content spreads. Watermark-only approaches fail interoperability because proprietary detection locks out third-party verification. Neither layer alone satisfies the regulation. Both are required.

The four criteria operate conjunctively. A system achieving excellent robustness but failing interoperability is non-compliant. This constraint eliminates certain architectural approaches entirely and narrows the design space for compliant systems.

### The four Article 50 criteria with practical examples

Each criterion addresses a specific aspect of AI content marking. All four must be satisfied for compliance.

| Effective | Interoperable |
|---|---|
| **Definition:**<br>Enables identification of AI-generated content in real-world use after distribution<br><br>✓ **Passes when:**<br>• Detection works after social media sharing<br>• Third parties can verify without special access<br><br>✗ **Fails when:**<br>• Only detectable at point of creation | **Definition:**<br>Works across different systems, platforms, and tools without proprietary dependencies<br><br>✓ **Passes when:**<br>• Open standard (C2PA) readable by any tool<br>• Public detection API available<br><br>✗ **Fails when:**<br>• Only provider's tools can verify (e.g., SynthID) |
| **Robust** | **Reliable** |
| **Definition:**<br>Survives typical content transformations encountered in distribution<br><br>✓ **Passes when:**<br>• Survives JPEG compression, cropping, screenshots<br>• Invisible watermarks persist through edits<br><br>✗ **Fails when:**<br>• Metadata stripped by screenshot (C2PA alone) | **Definition:**<br>Produces consistent, accurate detection results with low error rates<br><br>✓ **Passes when:**<br>• Low false positive rate (<0.1%)<br>• Consistent results across content types<br><br>✗ **Fails when:**<br>• Frequently flags human art as AI-generated |

*All four criteria must be met simultaneously. Failure on any single criterion = non-compliance.*

# Where providers stand today

Our assessment evaluated ten major AI providers against the four Article 50 criteria. The findings reveal a consistent pattern: partial compliance at best, with critical gaps in robustness or detection access.

**Metadata-first providers face robustness gaps.** OpenAI and Microsoft have implemented C2PA, the open standard for content provenance. This satisfies interoperability—any tool supporting C2PA can read the metadata—but fails robustness. Our testing confirms what providers acknowledge: C2PA metadata does not survive screenshots, most social media pipelines, or common editing workflows.

**Watermark-first providers face interoperability gaps.** Google's SynthID represents the most robust deployed watermarking system, having marked more than ten billion images. It survives compression, cropping, and screenshots. But detection requires Google's proprietary tools.

**Hybrid approaches come closest but still fall short.** Meta combines IPTC metadata with invisible watermarking—the right architectural direction. But detection remains platform-restricted rather than publicly accessible.

EXHIBIT 2

## No major provider is fully compliant today

### No major provider is fully compliant today

Assessment of ten providers against the four Article 50 criteria. No provider achieves compliance across all four.

● Compliant  ● Partial  ● Non-compliant

| Provider | Approach | Effective | Interoperable | Robust | Reliable |
|---|---|---|---|---|---|
| Google | SynthID watermark (proprietary) | ● | ● | ● | ● |
| OpenAI | C2PA metadata (DALL-E, GPT-4o) | ● | ● | ● | ● |
| Meta | Hybrid: IPTC + invisible watermark | ● | ● | ● | ● |
| Microsoft | C2PA metadata (Bing Image Creator) | ● | ● | ● | ● |
| Adobe | C2PA (Firefly) + Content Credentials | ● | ● | ● | ● |
| Stability AI | Invisible watermark + metadata | ● | ● | ● | ● |
| Midjourney | Metadata (steganographic) | ● | ● | ● | ● |
| Anthropic | Text only (no image gen) | ● | ● | ● | ● |
| xAI (Grok) | Visible watermark only | ● | ● | ● | ● |
| Kuaishou (Kling) | Visible watermark only | ● | ● | ● | ● |

Source: Provider documentation, public disclosures, and limited empirical testing. Assessment as of February 2026.

*Source: Provider documentation, public disclosures, and limited empirical testing. Assessment as of February 2026.*

# Ten principles for compliant architecture

Our analysis of regulatory requirements, technical specifications, and deployment experience at scale yields ten architectural principles for compliant marking systems.

## 1. Defense in depth

Compliant systems must implement at least three complementary marking layers, each covering the failure modes of others. The first layer—metadata such as C2PA—provides interoperability and rich provenance
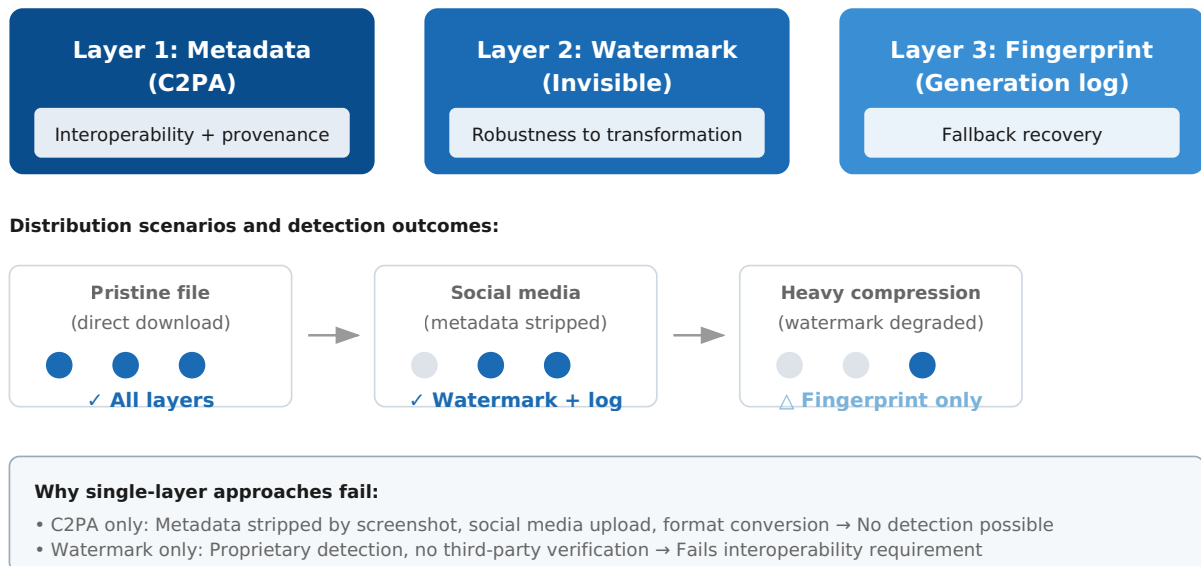
information. The second layer—an invisible watermark—provides robustness. The third layer—a fingerprint or generation log—enables fallback recovery when content is heavily degraded.

EXHIBIT 3

## Single-layer marking is explicitly non-compliant

### Single-layer marking is explicitly non-compliant

The Draft Code requires multi-layer marking. Each layer covers the failure modes of others.

| Layer 1: Metadata (C2PA) | Layer 2: Watermark (Invisible) | Layer 3: Fingerprint (Generation log) |
|---|---|---|
| Interoperability + provenance | Robustness to transformation | Fallback recovery |

**Distribution scenarios and detection outcomes:**

| Pristine file (direct download) | Social media (metadata stripped) | Heavy compression (watermark degraded) |
|---|---|---|
| ✓ All layers | ✓ Watermark + log | △ Fingerprint only |

**Why single-layer approaches fail:**
- C2PA only: Metadata stripped by screenshot, social media upload, format conversion → No detection possible
- Watermark only: Proprietary detection, no third-party verification → Fails interoperability requirement

## 2. Independent layer verification

Each marking layer provides independent evidence of AI generation. Detection systems should evaluate all available layers and report which were found.

## 3. Open verification

Third parties must be able to verify content without proprietary dependencies. The Draft Code's requirement for "free-of-charge interfaces" means detection cannot be gated behind API keys, bilateral agreements, or platform membership.

## 4. Marking by default

Marking must be architecturally enforced, not policy enforced. Systems offering opt-out create compliance gaps—bad actors will opt out.

## 5. Standards-based interoperability

Proprietary watermarks are acceptable for robustness, but interoperability requires a common layer that third parties can read without bilateral agreements. C2PA v2.3 represents the current industry convergence point.

## 6. Security through depth

Watermark security should not depend solely on secret algorithms, but detection interfaces must be protected against adversarial exploitation.

## 7. Ecosystem coordination

Marking fails if any link in the distribution chain breaks it. Providers mark at generation; platforms must preserve marks through their pipelines; deployers must surface labels to users.

## 8. Architectural flexibility

Today's standards will evolve. Compliant architectures should isolate standard-specific components for replacement without system redesign.

## 9. Privacy by design

The Draft Code is silent on privacy, but implementations must comply with GDPR. Privacy by design means no PII in default manifests, user control over identity disclosure, and a pathway for anonymous creation.

## 10. Cognitive simplicity

User-facing labels should minimize cognitive burden. A single icon indicating AI involvement, with detailed provenance accessible on demand, serves users better than complex taxonomies.

---

# Robustness requirements in practice

Watermarks must survive the transformations content encounters in real-world distribution. Our guidance draws on deployment experience at internet scale.

## The testing framework

Robustness evaluation requires discipline. Detection thresholds must be calibrated once and applied consistently across all transformations. Testing must cover both random variations within realistic ranges and worst-case extremes. And critically, testing must include transformation chains: watermarks surviving individual operations may fail when crop, scale, and compress are combined.

EXHIBIT 4

# Robustness testing requires 30+ transformations across eight categories

Comprehensive testing must include individual transforms, combination chains, and platform-specific pipelines.

| 30+ | 8 | 2-3 | 1 |
|---|---|---|---|
| Individual transforms | Categories | Chain depth to test | Operating point *(same threshold for all)* |

### Spatial
- Rotation (±5°, 90°, 180°, 270°)
- Flips (horizontal, vertical)
- Cropping (25%-90% retain)
- Scaling (50%-200%)

### Quality
- JPEG (Q30-Q90)
- WebP compression
- Format conversion
- Resolution change

### Color
- Brightness (±20%)
- Contrast adjustment
- Saturation/hue
- Grayscale conversion

### Noise
- Gaussian noise
- Salt & pepper
- Speckle, shot noise
- Impulse noise

### Filtering
- Gaussian blur (σ 0.5-2.0)
- Sharpening
- Denoising
- ↑ Often targets watermarks

### Overlay
- Text overlay
- Emoji overlay
- Logo/watermark overlay
- Partial occlusion

### Combination
- Crop + scale + compress
- Resize + noise + JPEG
- 2-3 transform chains
- **Often the binding constraint**

### Platform
- Instagram pipeline
- Twitter/X pipeline
- TikTok pipeline
- WhatsApp compression

■ Critical priority  ■ High priority  ■ Medium priority

*Note: Standard datasets underrepresent corner-case content. Explicitly test black-and-white images, logos, pixel art, gradients, and spars*

## Performance targets

Based on deployed systems, reasonable targets for true positive rates are: near-100% for identity (no transformation), 90-99% for individual transformations at random strengths, 75-95% for individual transformations at worst-case extremes, and 85-95% for combination chains and platform pipelines.

# Implementation checklist

## For AI Providers

### Marking Infrastructure

- ☐ Multi-layer marking implemented (metadata + watermark + fingerprint)
- ☐ C2PA manifests generated for all output
- ☐ Watermarks applied at generation, not post-processing
- ☐ Soft bindings registered for durable credentials
- ☐ No watermark-removal option in any interface

### Robustness

- [ ] Watermarks survive JPEG Q50 compression
- [ ] Watermarks survive screenshot capture
- [ ] Watermarks survive 25% crop
- [ ] Robustness validated through testing across all eight categories
- [ ] Combination attacks tested

### Detection Access

- [ ] Public verification interface available without authentication
- [ ] API documentation published
- [ ] Rate limits sufficient for platform-scale verification

### Privacy

- [ ] No PII in default manifests
- [ ] Privacy impact assessment completed
- [ ] Anonymous creation pathway exists
- [ ] GDPR alignment documented

## For Platforms

### Ingest

- [ ] C2PA metadata extracted on upload
- [ ] Watermark detection performed
- [ ] Soft-binding lookup attempted when metadata absent
- [ ] AI content flagged in internal systems

### Preservation

- [ ] Metadata preserved through processing
- [ ] Re-embedding attempted after lossy operations
- [ ] Provenance chain maintained for derivatives

### Disclosure

- [ ] AI-generated content labeled consistently
- [ ] Detailed provenance accessible on demand

---

# Timeline and risk

Enforcement begins August 2, 2026. Working backward from that date, the implementation windows are tight.

Adding C2PA metadata to existing generation pipelines typically requires two to four months. Implementing watermarking from scratch requires six to twelve months. Building detection APIs requires two to four months. Privacy review requires one to two months and should run in parallel.

> **Critical risk for providers without watermarking**
>
> Providers without watermarking infrastructure face the highest risk. The twelve-month upper bound for watermark implementation exceeds the time remaining. These organizations should begin immediately—or plan for interim compliance with metadata plus placeholder soft bindings, adding watermarks when ready.

# Two pathways to compliance

Organizations face a choice between minimal and best-practice compliance.

**EXHIBIT 5**

## Two pathways to compliance: minimal versus best practice



**Two pathways to compliance: minimal versus best practice**

Organizations must choose between baseline requirements and comprehensive implementation that creates margin

| Minimal Compliance | Best Practice |
|---|---|
| **Marking layers**<br>≥2 techniques (metadata + watermark) | **Marking layers**<br>≥3 layers with defined roles (metadata + watermark + fingerprint/log) |
| **Robustness testing**<br>Basic: compression, cropping, format conversion<br>No specific survival thresholds defined | **Robustness testing**<br>30+ transforms across 8 categories<br>Combination chains + platform pipelines |
| **Detection interface**<br>"Some form" of free access<br>Format and rate limits undefined | **Detection interface**<br>Public API with documented format<br>Sufficient rate limits for platform-scale use |
| **Privacy**<br>No explicit requirements in Draft Code | **Privacy**<br>Privacy by default, GDPR alignment documented |
| ⚠ **Risk exposure** Vulnerable to stricter enforcement interpretation | ✓ **Creates margin** Buffer against interpretation + future tightening |

*For most providers, the additional effort for best-practice compliance is warranted given regulatory uncertainty.*

For most providers, the additional effort for best-practice compliance is warranted given regulatory uncertainty.

**Minimal compliance** interprets requirements strictly from regulatory text: at least two marking techniques, survival through compression, cropping, and format conversion (no specific thresholds), some form of free

detection interface (format undefined). This approach is defensible against the literal text but carries risk if enforcement interprets requirements more broadly.

**Best-practice compliance** incorporates deployment experience and technical analysis: at least three marking layers with defined roles, comprehensive robustness testing across thirty-plus transformations including combinations, public API with documented format and rate limits, privacy by default with GDPR alignment, C2PA v2.3 with soft bindings.

---

# Looking ahead

The EU AI Act's transparency requirements mark a turning point for how AI-generated content must be handled. For providers, the mandate is clear: implement multi-layer marking, enable open verification, address privacy, and do it before August 2026.

The building blocks exist. C2PA v2.3 provides the interoperability layer. Watermarking at scale has been demonstrated. The soft-binding architecture bridges the gap between robustness and interoperability. What remains is implementation.

Organizations that act now will have time to iterate, test, and refine. Those that wait risk implementation under deadline pressure with no margin for error.

---

### About This Analysis

This guide synthesizes regulatory interpretation, provider landscape analysis, and architectural recommendations for EU AI Act Article 50 compliance. Sources include the December 2025 Draft Code of Practice, C2PA Technical Specification v2.3, WITNESS's privacy analysis, and Google DeepMind's published deployment experience.

The analysis focuses on image and video marking. Audio receives limited coverage. Text watermarking is acknowledged as technically immature and largely out of scope.

*This represents independent technical assessment and does not constitute legal advice.*

### References

1. Regulation (EU) 2024/1689 (EU AI Act), Article 50
2. European Commission, First Draft Code of Practice on Transparency of AI-Generated Content, December 2025
3. C2PA Technical Specification v2.3, December 2025
4. WITNESS, Privacy-First Transparency: Response to the First Draft EU AI Act Code of Practice, February 2026
5. Google DeepMind, SynthID-Image: Image watermarking at internet scale, arXiv, October 2025