

# How Does the EU's Digital Services Act Regulate Content Moderation? And Will it Work?

- Dr. Martin Husovec

# YouTube & copyright infringers



**Copyright Status**

Careful! You have 2 copyright strikes. One more and we may have to disable your account! [See details](#)

Type	Event	Content
Video	© STRIKE 1	<a href="#">LIVING IN A LOCKER SOTW - Viper Reforged #1   MINECRAFT HCF</a> Reason: Copyright takedown request Removed by ViperMC on Jan 1, 2019 Expires on Apr 1, 2019
	© STRIKE 2	<a href="#">WE ALLIED ON IMAKEMCVIDS - ViperHCF Reforged   Minecraft HCF</a> Reason: Copyright takedown request Removed by Jewell B Andrew on Jan 29, 2019 Expires on Apr 29, 2019

We may have to disable your account if you have 3 or more strikes within a 3 month period. [Learn more](#)

# Facebook & breastfeeding



# Twitter & Trump



**Donald J. Trump** ✓

@realDonaldTrump

45th President of the United States of

[Traducir la biografía](#)

📍 Washington, DC [Vote.DonaldJ](#)

49 Siguiendo 25,2 M Seguidores

*“The 75,000,000 great American Patriots who voted for me, AMERICA FIRST, and MAKE AMERICA GREAT AGAIN, will have a GIANT VOICE long into the future. They will not be disrespected or treated unfairly in any way, shape or form!!!”*

Shortly thereafter, the President Tweeted:

*“To all of those who have asked, I will not be going to the Inauguration on January 20th.”*

# AWS & hate speech

TECH • AMAZON

## Amazon Will Suspend Hosting For Pro-Trump Social Network Parler

Amazon's suspension of Parler's account means that unless it can find another host, once the ban takes effect on Sunday Parler will go offline.



**John Paczkowski**  
Technology and Business Editor



**Ryan Mac**  
BuzzFeed News Reporter

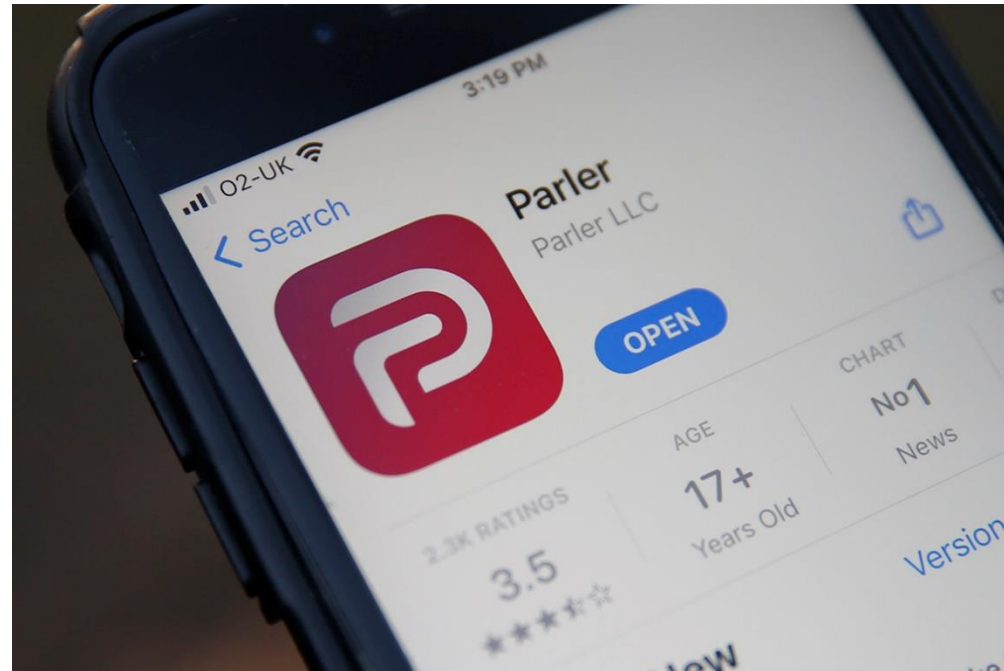
Updated on January 10, 2021 at 3:08 am

Posted on January 10, 2021 at 2:07 am

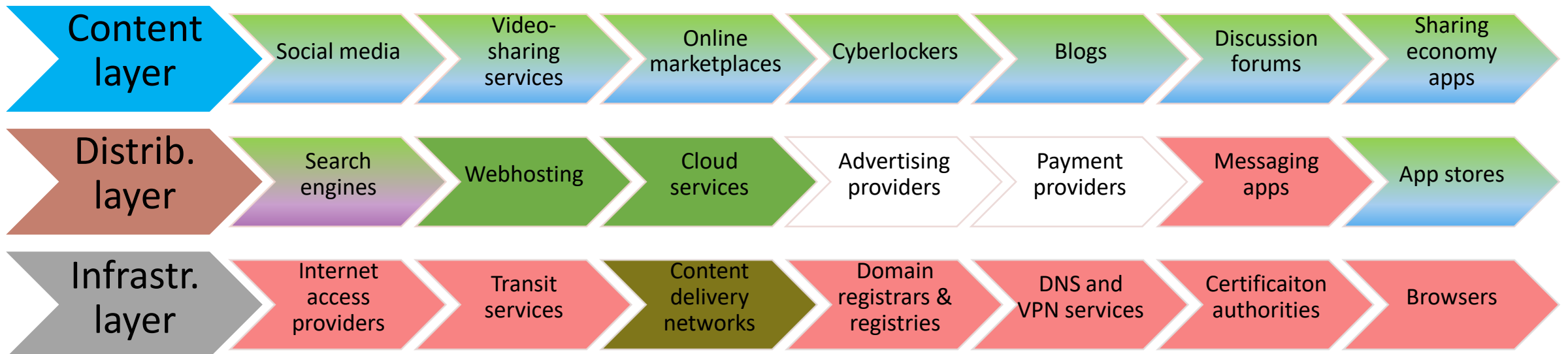


View 126 comments

# Apple App Store



# The main focus of regulation



Hosting, Online Platform, Mere Conduit, Caching, VLOSE

# Digital Services Act – from Feb 2024



Document 32022R2065



**Regulation (EU) 2022/2065 of the European Parliament and of the Council of 19 October 2022 on a Single Market For Digital Services and amending Directive 2000/31/EC (Digital Services Act) (Text with EEA relevance)**

PE/30/2022/REV/1

*OJ L 277, 27.10.2022, p. 1–102 (BG, ES, CS, DA, DE, ET, EL, EN, FR, GA, HR, IT, LV, LT, HU, MT, NL, PL, PT, RO, SK, SL, FI, SV)*

 In force

ELI: <http://data.europa.eu/eli/reg/2022/2065/oj>



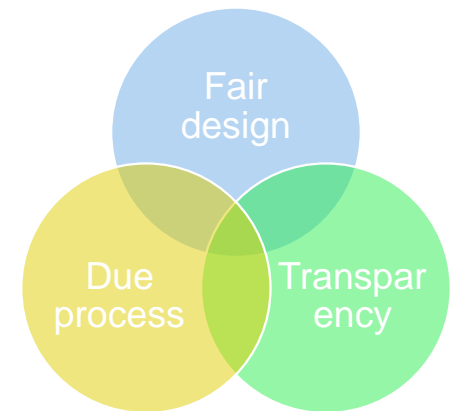
## A new generation of rules

	United States	European Union
I. generation: 1996-2000	Sec 230 CDA; Sec 512 DMCA	Articles 12-15 ECD
II. generation: 2020-?	? [PACT]	<b>Digital Services Act</b>

- **I. generation:** breathing space for speech & industries
  - Liability exemptions to avoid strict (or any) liability
- **II. generation:** regulation of risks posed by services
  - Regulatory expectations that overlay the liability social contract

# Due diligence obligations vs liability exemptions

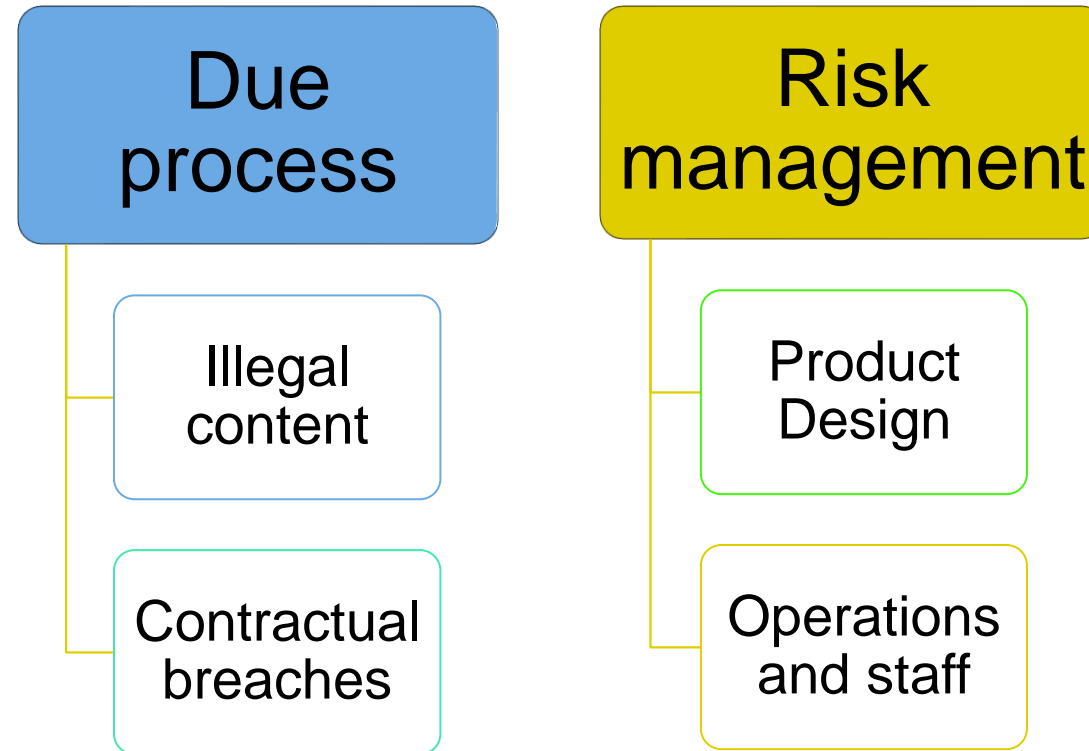
- The shift from **liability for content** to (legal) **accountability for the design** of services
  - Failing to comply does not lead to the illegality of service or liability for users' actions but targeted non-compliance with the stand-alone obligations
  - Even services which are *not* protected anymore by exemptions remain subject to due diligence obligations

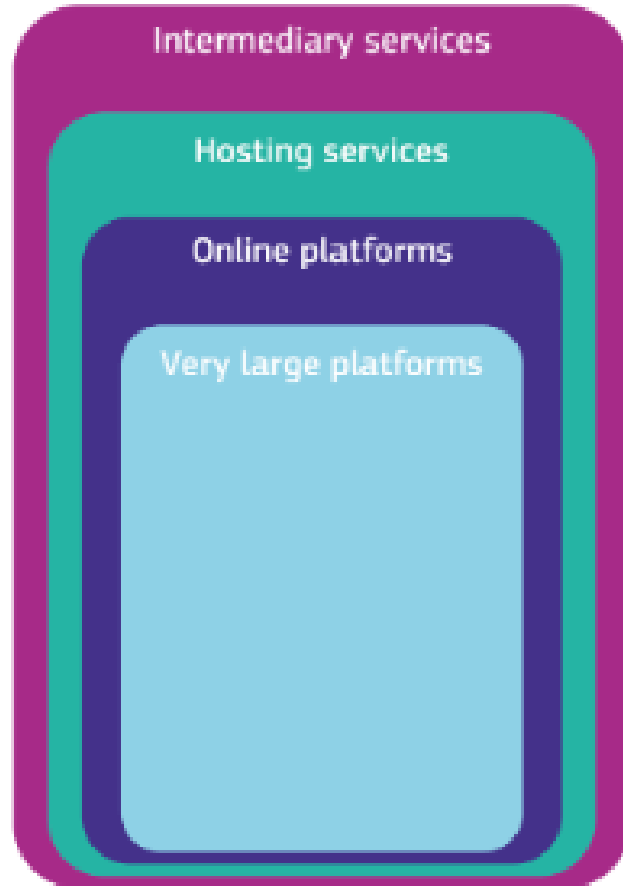


# The guts of the DSA

Obligations	Universal <i>All providers of conduit, caching, hosting services</i>	Basic <i>all hosting services</i>	Advanced <i>medium-to-large<sup>1</sup> online platforms</i>	Special <i>VLOPs &amp; VLOSEs</i>
<b>Content Moderation</b>	<b>Art 14</b> (fair content moderation)	<b>Art 16</b> (notice)  <b>Art 17</b> (statement of reasons)	<b>Art 20</b> (internal redress); <b>Art 21</b> (out-of-court mechanism); <b>Art 22</b> (trusted flaggers); <b>Art 23</b> (anti-abuse provisions); <b>Art 30-32</b> (specific rules on B2C marketplaces)	<b>Art 34-35</b> (risk mitigation assessment)  <b>Art 36</b> (crisis response mechanism)
<b>Fair Design</b> (user interfaces, recommender systems, advertising and other parts)	<b>Art 14</b> (fair content moderation)	<b>Art 16</b> (user-friendly notice and action)	<b>Art 25</b> (fair design of user-experience); <b>Art 26(3)</b> (advertising); <b>Art 27</b> (recommender systems); <b>Art 28</b> (protection of minors); <b>Art 30</b> (traceability of traders); <b>Art 31</b> (facilitating design for traders)	<b>Art 38</b> (recommender systems)  <b>Art 39</b> (risk mitigation assessment)
<b>Transparency</b>	<b>Art 15</b> (annual reporting)	<b>Art 17(5)</b> (database of all the statements of reasons)	<b>Art 22</b> (reports by trusted flaggers); <b>Art 24</b> (content moderation reports); <b>Art 26</b> (advertising disclosure)	<b>Art 39</b> (advertising archives); <b>Art 42</b> (content moderation transparency)
<b>Oversight</b>	<b>Art 11</b> (regulator's contact point); <b>Art 12</b> (recipient's contact point); <b>Art 13</b> (legal representative)	<b>Art 18</b> (notification of suspected relevant crimes)	(-)	<b>Art 37</b> (auditing); <b>Art 40</b> (data access/scrutiny); <b>Art 41</b> (compliance function)

## DSA's two main tools:





Overview of the DSA's due diligence obligations

Obligations	Universal	Basic	Advanced	Special
	All providers of <u>conduit</u> , <u>caching</u> , <u>hosting services</u>	all hosting services	medium-to-large <sup>1</sup> online platforms	VLOPs & VLOSEs

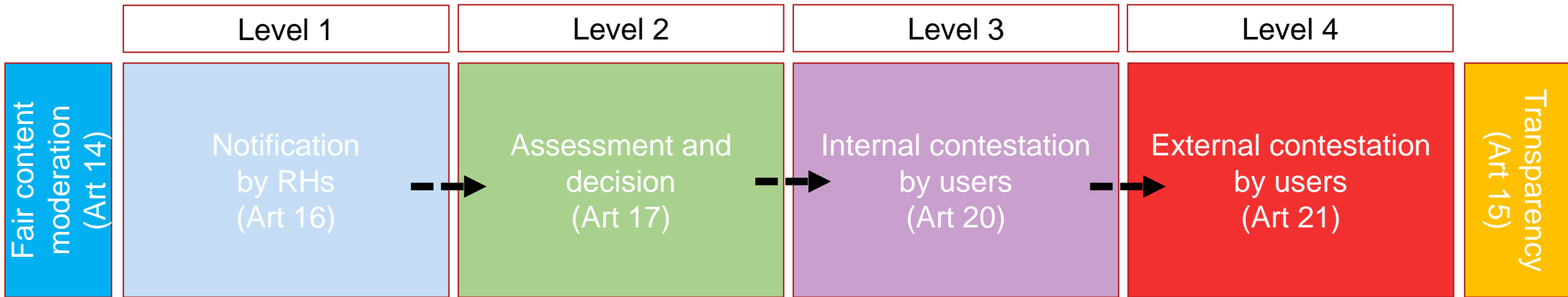
- **Online platforms** = medium-to-large firms (50+ employees or turnover 10+ million EUR)
- **VLOPs/VLOSEs** = 45+ mil average active monthly EU users

# VLOPs / VLOSEs

	Company	Digital Service	Type	Est. (cc)	Users (mil)	User-generated-content components
<b>Search</b>	Alphabet <sup>11</sup>	Google Search	VLOSE	IE	332+	Paid and unpaid search results
	Microsoft <sup>12</sup>	Bing	VLOSE	IE	107	Paid and unpaid search results
<b>Social media</b>	Alphabet	YouTube	VLOP	IE	401+	Videos, sound, photos & text
	Meta <sup>13</sup>	Facebook	VLOP	IE	255	Videos, sound, photos & text
	Meta	Instagram	VLOP	IE	250	Videos, sound, photos & text
	Bytedance <sup>14</sup>	TikTok	VLOP	IE	125	Videos, sound, photos & text
	Microsoft	LinkedIn	VLOP	IE	122	Videos, sound, photos & text
	Snap <sup>15</sup>	Snapchat	VLOP	?	96+	Videos, sound, photos & text
	Pinterest <sup>16</sup>	Pinterest	VLOP	?	n/a	Videos, sound, photos & text
	Twitter <sup>17</sup>	Twitter	VLOP	?	100+	Videos, sound, photos & text
<b>App stores</b>	Alphabet	Google App Store	VLOP	IE	274+	Mobile apps
	Apple <sup>18</sup>	Apple App Store	VLOP	IE	n/a	Mobile apps
<b>Wiki</b>	Wikimedia <sup>19</sup>	Wikipedia	VLOP	?	151+	Mostly text and photos
<b>Markets</b>	Amazon <sup>20</sup>	Amazon Marketplace	VLOP	LX	n/a	Sellers' offerings & users' reviews
	Alphabet	Google Shopping	VLOP	IE	74+	Sellers' offerings & users' reviews
	Alibaba <sup>21</sup>	AliExpress	VLOP	?	n/a	Sellers' offerings & users' reviews
	Booking.com <sup>22</sup>	Booking.com	VLOP	NL	n/a	Sellers' offerings & users' reviews
<b>Maps</b>	Alphabet	Google Maps	VLOP	IE	278+	Shop profiles, reviews, etc.

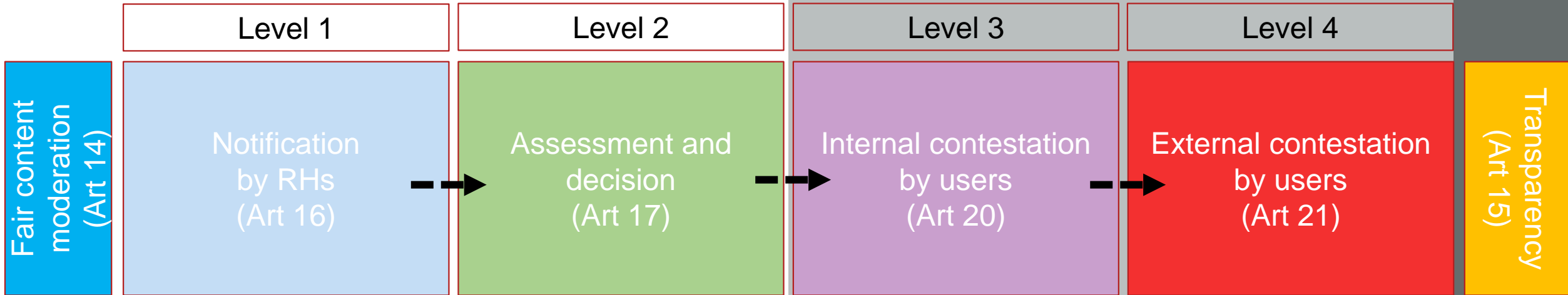
# Content Moderation

# Content Moderation





# Content Moderation



**Hosting services  
& any-size**

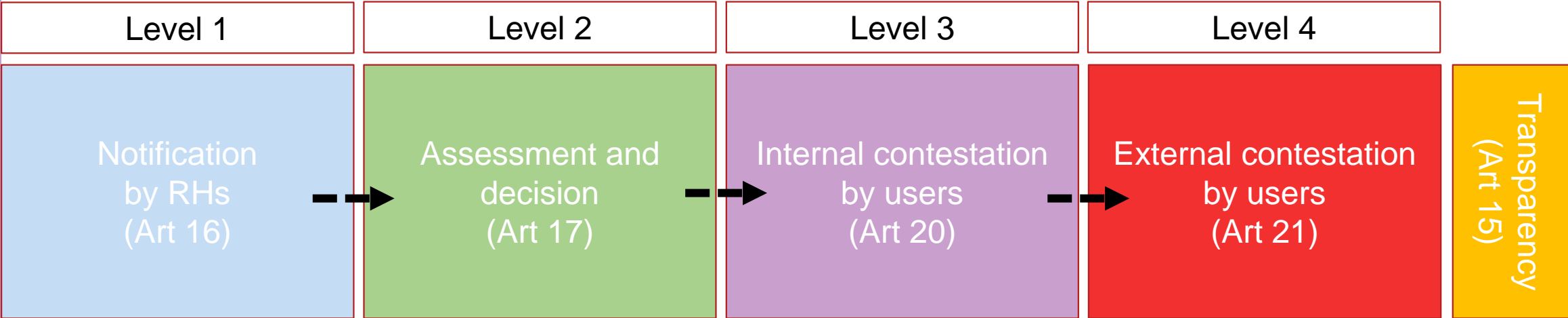
**Online platforms  
& mid-size +**

**mid-size +**

# Content Moderation

Rule  
formation

Fair content  
moderation  
(Art 14)



## Two main obligations

- Codification & explanation of all restrictions
  - “any restrictions that they impose in relation to the use of their service in respect of” UGC content
- Conduct content moderation fairly
  - act “diligently, objectively and proportionately” with due regard to the fundamental rights of others

## Article 14(1)

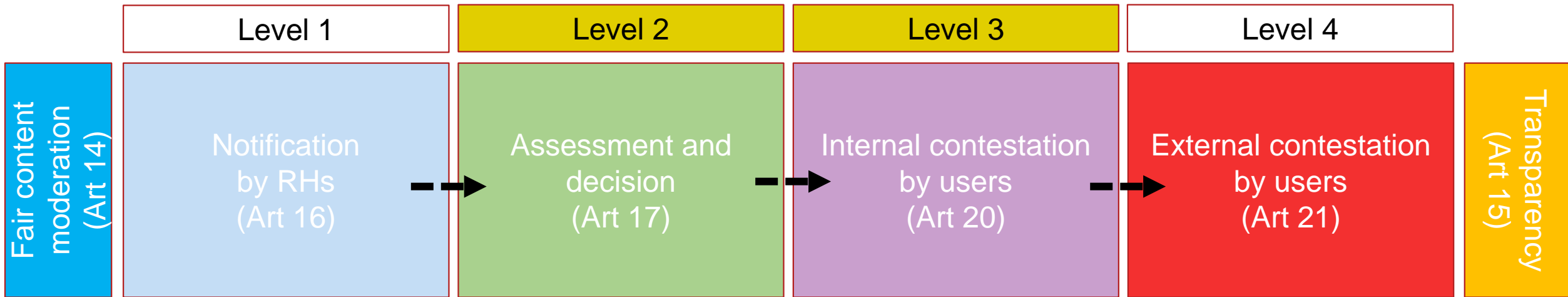
“Providers of intermediary services shall **include information on any restrictions that they impose in relation to the use of their service in respect of information provided by the recipients of the service, in their terms and conditions.** That information shall **include** information on any policies, procedures, measures and tools used **for the purpose of content moderation**, including algorithmic decision-making and human review, as well as the rules of procedure of their internal complaint handling system. **It shall be set out in clear, plain, intelligible, user-friendly and unambiguous language, and shall be publicly available in an easily accessible and machine-readable format.**”

## Article 14(4) & Recital 47

“Providers of intermediary services shall act in a diligent, objective and proportionate manner in applying and enforcing the restrictions referred to in paragraph 1, with due regard to the rights and legitimate interests of all parties involved, including the fundamental rights of the recipients of the service”

(47) When designing, applying and enforcing those restrictions  
(..).

# Content Moderation



## Level 2 & 3 (internal operations)

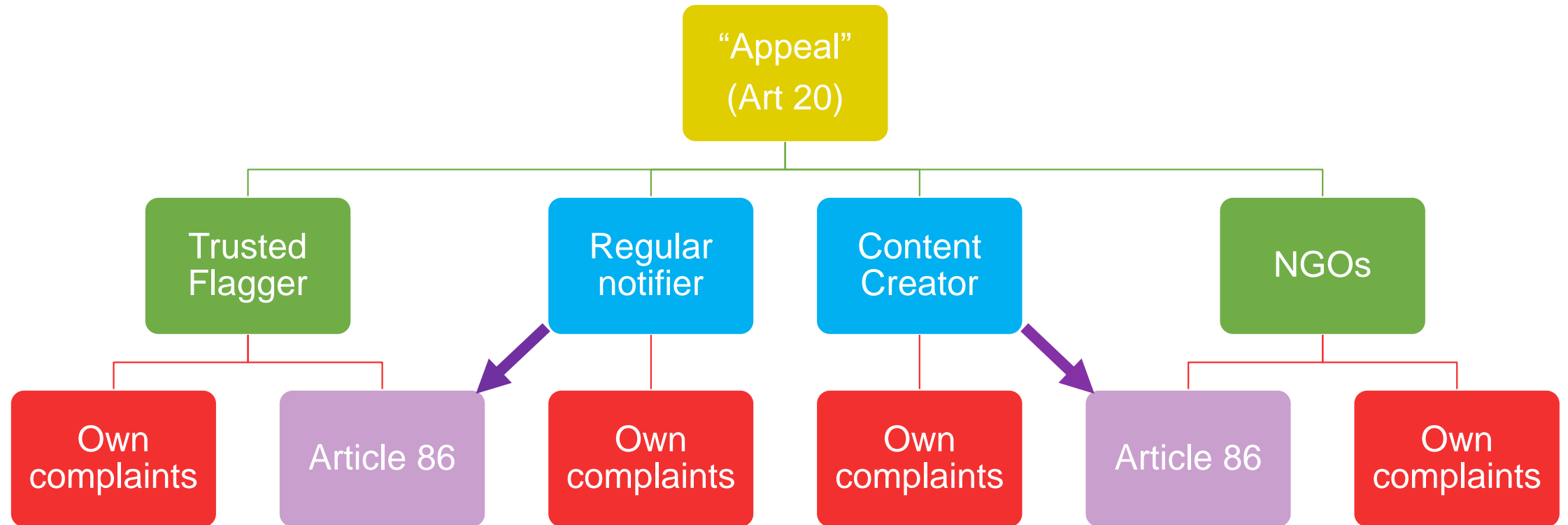
- Initial decision (Art 17)
  - providers must issue a statement of reasons
  - broad notion of COMO: visibility, monetization, etc.
  - specific explanation of reasons
  - can be automated (but can small companies do this without vendors?)
- Internal appeal mechanism (Art 20)
  - “easy to access and user-friendly”
  - “not solely automated” + timely, diligent and objective manner
  - inform affected parties

## A broad scope of relevant restrictions

<b>Content visibility restrictions (17(1)(a))</b>	<b>Monetisation restrictions (17(1)(b))</b>	<b>Service restriction (17(1)(c))</b>	<b>Account restriction (17(1)(d))</b>	<b>Content moderation outside Article 17</b>
Removing content	Forfeiting earnings	Blocking access to a service	Blocking accounts	Shaming
Suspending content	Suspending access to earnings	Suspending rights on a service	Suspending accounts	Community service
Relocating content	Disabling income on some content	Reducing speed on a service	Erasing accounts	Attach context
Redacting content	Reducing earnings	Limits puts on notifications	Blacklisting registration	Attach warning
Shadow banning or Blacklisting	Fining	Limits on internal complaints		Apology, labelling



# Who can complain?



## Transparency for mid-sized companies

- Annual reports (or bi-annual for VLOPs) about the content moderation practices (Article 15)
- Aggregate Lumen-like database for all statements of reasons (Article 24(5)) for platforms
  - Will this work? How to anonymise this? But the structuring needed to comply with this can be helpful for standardisation, including annual reports, data access, and risk management

## External review: ADR bodies

- Follows broadly Fiala & Husovec 72 (2022) *International Review of Law and Economics* (re funding)
- The national regulators certify ADR providers
  - Up to market or states to create them
- Content creators & notifiers can complain
  - IF they win, provider reimburses the fee + costs
  - IF they lose, they pay the fee

Level 4

External contestation  
by users  
(Art 21)

## Article 21

- ADR (Art 21): out-of-court settlement bodies
  - Regulators certify entities (must be independent of Ps & users)
    - FB's Oversight Board is clearly not independent in this sense
  - Content creators & notifiers (and their reps) can use the option
  - ADR provider is complainants' choice; no need to exhaust appeals
  - ADR issue decisions: non-binding, Ps must engage in good faith
  - P compensates complainants who win (pays fees & possibly costs)
  - Complainants that lose pays their own fees & costs

## The effect of ADR

- External interpreter of platform's rules (= loss of power)
- Incentive for platforms to be clearer (= push to codify)
- Incentive for platforms to resolve internally (= costs)

### **BUT:**

- Does not take away the power to make rules!

## Rule-making vs Interpretation

- The basic rule: the **proceduralist approach** constrains interpretation but not rule-making of platforms
- DSA takes mostly proceduralist approach (and a systems-design approach for risk management), with the exception of Art 14(4)
- But my view: Article 14(4): constrains only **arbitrary** & **grossly disproportionate rule-making**; all other legal rule-making is fine
- However, once a content rule is expressed by Ps, they don't decide its meaning unless they change it again (as we do with regular contracts)

## Case 1: VIP No-Moderation List

- A micro-blogging site decides to create a list of VIPs whose content is not moderated at all (regardless of illegality or contractual nature)
- VIPs are all top elected officials in all countries of the UN
- Is the policy in violation of Article 14(1)?
  - *IMO: not, if properly described.*
- Is the policy in violation of Article 14(4)?
  - *IMO: yes, due to the impact of illegal content (separate from Art 6!)*

## Case 2: Pay2Say

- A micro-blogging site has a new Pay2Say product
- For 5 EUR a month, you can say whatever you want on the service, as long it is legal in your country (= no contractual restriction on speech [e.g., disinformation, nudity, vulgar content], only illegality).
- Everyone else is moderated on ToS violations & illegality.
- Is the policy in violation of Article 14(1)? *No if disclosed.*
- Is the policy in violation of Article 14(4)? *Probably no.*



# Open issues 1

- Statement of reasons
  - Automation by small providers (part of licensed COMO solutions?)
  - Recommended system changes vs statement of reasons (individualised)
  - FB page owners (eg news orgs) with their own COMO as hosting services
  - Transparency reporting and its standardisation & real-time transmission
- ADR
  - Certification of fee structures, calculation of reasonable costs
  - Transparency obligations, oversight & abuse
  - Scope of specialisation by ADRs

## Open issues 2

- Rule-making:
  - Article 14(4) – what is disproportionate?
  - Article 14(1) – technical constrains?
  - Article 14(1) – what is “clear, plain, intelligible, user-friendly and unambiguous language” vs doable
    - Trade-off between: administrability (scalability) & explainability
- Private Enforcement of COMO due diligence obligations
  - Impact on contract law
  - Impact on private claims, e.g., copyright holders

## Risk management and its impact on COMO

- VLOPs / VLOSEs are subject to additional requirements
  - mostly extended or intensified reporting obligations
  - researchers' data access
  - unique obligations: profiling-free choice on recommender systems, or advertising archives
- **MAIN: a regulatory dialogue about risk management**

# VLOP

	Company	Digital Service	Type	Est. (cc)	Users (mil)	User-generated-content components
<b>Search</b>	Alphabet <sup>11</sup>	Google Search	VLOSE	IE	332+	Paid and unpaid search results
	Microsoft <sup>12</sup>	Bing	VLOSE	IE	107	Paid and unpaid search results
<b>Social media</b>	Alphabet	YouTube	VLOP	IE	401+	Videos, sound, photos & text
	Meta <sup>13</sup>	Facebook	VLOP	IE	255	Videos, sound, photos & text
	Meta	Instagram	VLOP	IE	250	Videos, sound, photos & text
	Bytedance <sup>14</sup>	TikTok	VLOP	IE	125	Videos, sound, photos & text
	Microsoft	LinkedIn	VLOP	IE	122	Videos, sound, photos & text
	Snap <sup>15</sup>	Snapchat	VLOP	?	96+	Videos, sound, photos & text
	Pinterest <sup>16</sup>	Pinterest	VLOP	?	n/a	Videos, sound, photos & text
<b>App stores</b>	Alphabet	Google App Store	VLOP	IE	274+	Mobile apps
	Apple <sup>18</sup>	Apple App Store	VLOP	IE	n/a	Mobile apps
<b>Wiki</b>	Wikimedia <sup>19</sup>	Wikipedia	VLOP	?	151+	Mostly text and photos
<b>Markets</b>	Amazon <sup>20</sup>	Amazon Marketplace	VLOP	LX	n/a	Sellers' offerings & users' reviews
	Alphabet	Google Shopping	VLOP	IE	74+	Sellers' offerings & users' reviews
	Alibaba <sup>21</sup>	AliExpress	VLOP	?	n/a	Sellers' offerings & users' reviews
	Booking.com <sup>22</sup>	Booking.com	VLOP	NL	n/a	Sellers' offerings & users' reviews
<b>Maps</b>	Alphabet	Google Maps	VLOP	IE	278+	Shop profiles, reviews, etc.

## Risk Management Dialogue

- Regulatory dialogue put in place due to the opacity of the ecosystem & information asymmetry
- The regulator has no clear idea of risks, or contributing factors, and is in dark about solutions
- Forces providers to think about this, let themselves be reviewed by others (auditors, researchers, field NGOs), and then the regulator forms an opinion



## VLOP's risk management: Article 34(1)

Providers of very large online platforms and of very large online search engines shall **diligently identify, analyse and assess any systemic risks in the Union** stemming from the design or functioning of their service and its related systems, including algorithmic systems, or from the use made of their services. This risk assessment shall be **specific** to their services and **proportionate** to the systemic risks, taking into consideration their **severity and probability**, and **shall include** the following systemic risks: (..)

## VLOP's risk mitigation

Risk Areas & Categories	Recommender systems	Content moderation	Terms and conditions	Advertising	Data practices	Other areas
<b>Illegal content</b> (Art 34(1)(a))	<i>Examples:</i> (a) terrorist content; (b) child sexual abuse; (c) illegal hate speech; (d) intellectual property infringements; (e) defamation; (f) sale of unsafe products; (g) cyberstalking or grooming; or (h) <i>any other</i> areas of illegal content or behaviour.					
<b>Fundamental rights</b> (Art 34(1)(b))	<i>Examples:</i> <sup>97</sup> (a) human dignity; (b) freedom of expression and information, including media freedom and pluralism; (c) right to private life; (d) data protection; (e) right to non-discrimination; (f) rights of the child; (g) consumer protection; or (h) <i>any other</i> fundamental rights.					
<b>Public security and elections</b> (Art 34(1)(c))	<i>Exhaustive subcategories:</i> <sup>98</sup> (a) civic discourse; (b) electoral process; and (c) public security.					
<b>Health and well-being</b> (Art 34(1)(d))	<i>Exhaustive subcategories:</i> <sup>99</sup> (a) gender-based violence; (b) public health; (c) rights of Minors; (d) physical well-being; and (e) mental well-being.					



## Metaphor: safety regulation of public protests

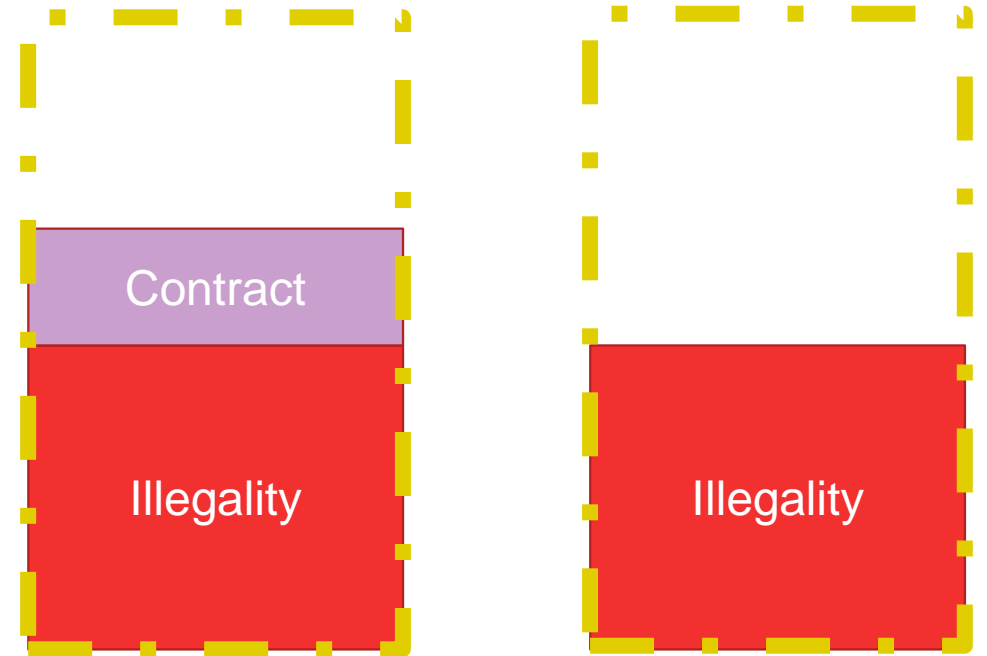
The authorities can (and do) partly restrict **how and when protest activities take place** (eg streets, hours, or use of amplification tools), and **take measures to prevent harm to protesters or others** (eg boost police presence), but cannot select speakers, or dictate the content





## Thus

- Two social media companies A & B set their baseline of rules [eg disinformation] and mitigate the risks on that basis
- Since parliaments permit both models of rules, they also permit two mixtures; risk mitigation cannot negate the existence of such choice



# Risk Mitigation Measures

Type of content	Priority by type of intervention	Examples
<b>1. Illegal content</b> (e.g., hate speech, terrorist content, copyright infringement)	1. Content removal	removal of content
	2. Visibility restrictions	age-gating of content or recommendations
	3. Nudges and incentives	demonetization of content or de-ranking borderline content
	4. Empowerment	flagging systems; rating systems; help lines; information
<b>2. Legal content</b> (e.g., disinformation, sensitive content, nude content, gambling content)	1. Empowerment	choice on recommendations; parental consent; rating systems; hotlines; information; suggestion tools;
	2. Nudges and incentives	default on recommendations; costly super-sharing; verification; pre-publication notices; parental consent; demonetization of content
	3. Visibility restrictions	age-gating of content or recommendations; de-ranking
	4. Content removal	removal of content



**Want to know more?**

[husovec.eu/DSA](https://husovec.eu/DSA)

[LSE Short Course on the EU Digital Services Act](#)