

ALEXANDER PIRANG

Freedom of Expression in the Platform Society

Internet und Gesellschaft



Mohr Siebeck

Internet und Gesellschaft
Schriften des Alexander von Humboldt Institut
für Internet und Gesellschaft

Herausgegeben von

Jeanette Hofmann, Matthias C. Kettemann,
Björn Scheuermann, Thomas Schildhauer
und Wolfgang Schulz

38



Alexander Pirang

Freedom of expression in the platform society

The right to freedom of expression as a constraint on
public authorities' power to govern online speech
through platform companies

Mohr Siebeck

Alexander Pirang, born 1989; law studies at Bucerius Law School, Hamburg and in Saint Petersburg; legal clerkship at the Berlin Court of Appeal; research associate at the Alexander von Humboldt Institute for Internet and Society; non-resident fellow with the Global Public Policy Institute; since 2021 policy officer at the Federal Ministry for Economic Affairs and Climate Action; 2023 doctorate (University of Hamburg).

Open Access gefördert durch den Fachinformationsdienst (FID) interdisziplinäre Rechtsforschung in Berlin. / Open Access funded by the Fachinformationsdienst (FID) interdisziplinäre Rechtsforschung in Berlin.

ISBN 978-3-16-163962-3 / eISBN 978-3-16-163963-0

DOI 10.1628/978-3-16-163963-0

ISSN 2199-0344 / eISSN 2569-4081 (Internet und Gesellschaft)

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliographie; detailed bibliographic data are available at <https://dnb.dnb.de>.

Published by Mohr Siebeck Tübingen, Germany, 2024. www.mohrsiebeck.com

© Alexander Pirang

This publication is licensed under the license “Creative Commons Attribution – ShareAlike 4.0 International” (CC BY-SA 4.0). A complete Version of the license text can be found at: <https://creativecommons.org/licenses/by-sa/4.0/>. Any use not covered by the above license is prohibited and illegal without the permission of the author.

The book was typeset using Times typeface and printed on non-aging paper by Laupp & Göbel in Gomaringen. It was bound by Nädle in Nehren.

Printed in Germany.

Vorwort (Acknowledgements)

Die vorliegende Arbeit wurde von der Fakultät für Rechtswissenschaft der Universität Hamburg im August 2023 als Dissertation angenommen. Rechtsprechung und Literatur konnten bis zur Einreichung im Januar 2023 umfassend berücksichtigt werden. Bis Februar 2024 veröffentlichte Literatur wurde zum Teil ergänzend aufgenommen.

Herzlich danken möchte ich zuallererst Professor Dr. Wolfgang Schulz, der mich mit wertvollen Denkanstößen wiederholt auf den richtigen Weg gebracht und mit seiner humorvollen und wertschätzenden Art für notwendige Motivationsschübe gesorgt hat. Einen besseren Doktorvater hätte ich mir nicht wünschen können.

Großer Dank gebührt auch Professor Dr. Hans-Heinrich Trute, der die Mühe des Zweitgutachtens auf sich genommen hat. Den Herausgeber:innen Prof. Dr. Jeanette Hofmann, Prof. Dr. Matthias C. Kettemann, Prof. Dr. Björn Scheuermann, Prof. Dr. Thomas Schildhauer und Prof. Dr. Wolfgang Schulz danke ich für die Aufnahme der Arbeit in die Schriftenreihe Internet und Gesellschaft. Dem Fachinformationsdienst für internationale und interdisziplinäre Rechtsforschung danke ich für die großzügige Förderung der Open-Access-Publikation.

Die Arbeit ist in substanziellen Teilen am Alexander von Humboldt Institut für Internet und Gesellschaft entstanden. Das interdisziplinäre und inspirierende Umfeld am Institut hat mich dazu ermutigt, in meiner Dissertation auch über den juristischen Tellerrand hinaus zu blicken.

Dass das mitunter einsame Unterfangen einer Doktorarbeit durch anregende Gespräche aufgelockert wurde, habe ich meinen Freund:innen zu verdanken. Thorsten Benner und Klaas Eller haben mich entscheidend darin bestärkt, das Projekt Doktorarbeit anzugehen. Tom Langerhans hat die Promotionsphase auf seine geradlinige und kluge Art am engsten begleitet.

Meinen Eltern danke ich für Ihre unbedingte Loyalität und ihren Glauben an mich.

Der größte Dank gebührt meiner Frau Anna, die mir auch nach der Geburt unserer Kinder die Freiräume, die eine zeitintensive wissenschaftliche Arbeit erfordert, ermöglicht hat. Ohne diese, keinesfalls selbstverständliche Unterstützung

und ihre wortreichen Ermahnungen, irgendwann zum Ende zu kommen, hätte die Arbeit nicht gelingen können.

Meine Töchter Aglaia und Flora führen mir tagtäglich vor Augen, wie wichtig es ist, mit Neugier und Interesse auf die Welt zu blicken. Ihnen ist diese Arbeit gewidmet.

Berlin im März 2024

Alexander Pirang

Chapter breakdown

Vorwort (Acknowledgements)	V
Table of contents	IX
List of figures and tables	XVII
List of abbreviations	XIX
Introduction	1
<i>I. Setting the stage</i>	1
<i>II. Hypothesis and research questions</i>	4
<i>III. Research scope</i>	5
<i>IV. Research Design</i>	12
<i>V. Conceptual clarification: governance and regulation</i>	20
Chapter 1: First-order online speech governance: content moderation	23
<i>I. What is content moderation?</i>	23
<i>II. What role does context play in content moderation?</i>	55
<i>III. To what extent does content moderation impact freedom of expression?</i>	60
<i>IV. How do different incentives and constraints influence platform companies' decision-making regarding content moderation?</i>	80
<i>V. Conclusion</i>	99

Chapter 2: Second-order online speech governance: platform regulation	105
<i>I. How has platform regulation evolved in the EU?</i>	105
<i>II. What are the rationales behind platform regulation?</i>	113
<i>III. How does platform regulation interact with other influences on platform companies' decision-making?</i>	116
<i>IV. How does platform regulation operate in practice?</i>	120
<i>V. Conclusion</i>	196
Chapter 3: Third-order online speech governance: freedom of expression	201
<i>I. To what extent does platform regulation allow public authorities to launder state action through platform companies' private ordering?</i>	202
<i>II. To what extent does the case law of the ECtHR and the CJEU on causation and attribution provide avenues for overcoming the risks of laundered state action in the context of platform regulation?</i>	223
<i>III. How can limitations of users' right to freedom of expression be established where public authorities require platform companies to moderate content?</i>	290
<i>IV. To what extent are public authorities obligated to minimize the risks of over-blocking in the context of platform regulation in order to justify limitations of users' right to freedom of expression?</i>	311
Conclusion	371
Deutsche Kurzzusammenfassung	379
Table of legislation	383
Table of cases	385
Literature	391
Index	429

Table of contents

Vorwort (Acknowledgements)	V
Chapter breakdown	IX
List of figures and tables	XVII
List of abbreviations	XIX
Introduction	1
<i>I. Setting the stage</i>	1
<i>II. Hypothesis and research questions</i>	4
<i>III. Research scope</i>	5
1. Content moderation by online platforms	5
a) Online platforms	5
b) Content moderation	7
c) Application layer	10
2. Platform regulation	10
3. Fundamental rights framework	11
<i>IV. Research Design</i>	12
1. General approach: three orders of online speech governance	13
2. First-order online speech governance: content moderation	14
3. Second-order online speech governance: platform regulation	15
4. Third-order online speech governance: right to freedom of expression	17
<i>V. Conceptual clarification: governance and regulation</i>	20
Chapter 1: First-order online speech governance: content moderation	23
<i>I. What is content moderation?</i>	23
1. Content moderation process	24
2. Policy development stage	26

a) Overview of platform law	27
b) Challenges of policy development at scale	28
c) Transparency	30
3. Enforcement stage	31
a) Content review and decision-making	32
aa) Human review	32
bb) Automated filters	33
(1) Matching	34
(2) Classifying	36
(3) Challenges of automated content moderation	38
(a) Accuracy	38
(b) Bias	40
(c) Explainability	43
cc) Decision-making	44
b) Sanctioning	46
aa) Content-level restrictions	46
bb) Account-level sanctions	47
cc) Visibility restrictions	48
dd) Other	50
c) Timing	51
d) Transparency	52
e) Redress	54
<i>II. What role does context play in content moderation?</i>	55
<i>III. To what extent does content moderation impact freedom of expression?</i>	60
1. The right to freedom of expression under the Convention and the Charter	61
2. Doctrinal challenges posed by private ordering	64
3. Restrictive effects of content moderation on users' freedom of expression	67
a) Scope of analysis	67
b) Content-level restrictions	68
aa) Hard control	68
bb) Soft control	70
c) Account-level restrictions	71
d) Visibility restrictions	72
aa) Search, viewing, and engagement constraints	72
bb) Algorithmic downranking	74
e) Chilling effects	79

<i>IV. How do different incentives and constraints influence platform companies' decision-making regarding content moderation?</i>	80
1. Conceptual starting points	81
2. Social norms	83
3. Economic incentives	88
4. Constraints resulting from limited accuracy of content moderation at scale	91
5. Towards a systematic understanding of non-regulatory influences on content moderation	93
<i>V. Conclusion</i>	99
Chapter 2: Second-order online speech governance: platform regulation	105
<i>I. How has platform regulation evolved in the EU?</i>	105
1. Starting point: horizontal safe harbor rules in the E-Commerce Directive	106
2. Transition period: sectoral regulation by the EU and member states	109
3. Back to the future: Digital Services Act	111
<i>II. What are the rationales behind platform regulation?</i>	113
<i>III. How does platform regulation interact with other influences on platform companies' decision-making?</i>	116
<i>IV. How does platform regulation operate in practice?</i>	120
1. Scope of analysis and terminology	121
2. Conditional liability	123
a) Regulatory strategy	123
b) Conditional liability in practice	124
aa) Digital Services Act: conditional liability regime	124
(1) Background and objectives	124
(2) Scope	125
(3) Conditional liability exemption of hosting services	126
bb) Directive on Copyright in the Digital Single Market	130
(1) Background and objectives	130
(2) Scope	132
(3) Primary liability for unauthorized communicating of works to the public	132
(4) Conditional liability exemption	134

(5) Redress	138
c) Assessment	139
3. Soft law and co-regulation	142
a) Regulatory strategy	142
b) Co-regulation in practice: EU Code of Conduct on Countering Illegal Hate Speech Online	145
aa) Background and objectives	145
bb) Scope	146
cc) Commitments regarding notice and action	147
dd) Monitoring and enforcement	148
c) Assessment	151
4. Command-and-control regulation	153
a) Regulatory strategy	153
b) Command-and-control regulation in practice	155
aa) German Network Enforcement Act: notice and action	155
(1) Background and objectives	155
(2) Scope	157
(3) Notice and action mechanism	157
(4) Transparency obligations	164
(5) Redress	167
(6) Monitoring and enforcement	168
bb) Digital Services Act: notice and action mechanism	170
(1) Background and objectives	170
(2) Scope	171
(3) Notice and action mechanism	171
(4) Transparency obligations	174
(5) Redress	176
(6) Monitoring and enforcement	178
cc) Digital Services Act: systemic risk mitigation	180
(1) Background and objectives	180
(2) Scope	180
(3) Systemic risk assessment and mitigation	181
(4) Transparency obligations	186
(5) Monitoring and enforcement	188
c) Assessment	190
<i>V. Conclusion</i>	196

Chapter 3: Third-order online speech governance: freedom of expression	201
<i>I. To what extent does platform regulation allow public authorities to launder state action through platform companies' private ordering?</i>	202
1. Interference criteria	203
2. Doctrinal challenges posed by platform regulation	205
a) Discretionary standards	206
b) Non-binding demands	209
c) Meta-regulation	210
d) Non-regulatory influences	213
e) Over-blocking	215
3. Causation and attribution as potential avenues for holding public authorities responsible	216
a) Causation and attribution in the context of online speech governance	216
b) Implications of procedural differences between the ECtHR and the CJEU	219
<i>II. To what extent does the case law of the ECtHR and the CJEU on causation and attribution provide avenues for overcoming the risks of laundered state action in the context of platform regulation?</i>	223
1. Direct responsibility of public authorities	223
a) Causation in public international law	224
b) ECtHR case law on causation	226
aa) Factual and legal causation in multi-actor scenarios	226
bb) Causation in internet-related case law	229
(1) Access blocking	229
(2) Intermediary liability	231
cc) Multiple causes and risk	236
c) CJEU case law on causation	239
aa) Factual and legal causation in multi-actor scenarios	239
bb) Causation in internet-related case law	241
d) Summary and discussion	248
aa) General principles	248
bb) Risk of laundered state action in light of the case law analysis	249
(1) Discretionary standards	249
(2) Non-binding demands	251
(3) Meta-regulation	252
(4) Non-regulatory influences	254
(5) Over-blocking	255

2. Vicarious responsibility of public authorities for private conduct	257
a) Attribution in public international law	257
aa) General attribution principles	258
bb) Attributable conduct of non-state entities under the ARSIWA	259
(1) Legal mandate to exercise governmental authority	259
(2) Instruction, Direction, or Control	261
b) ECtHR case law on attributable conduct of non-state entities	263
aa) Privatization of public functions	265
bb) Increasing focus on vicarious state responsibility	268
(1) The Radio France criteria	268
(2) The Kotov criteria	272
c) CJEU case law on attributable conduct of non-state entities	276
aa) Attribution of private conduct in the context of fundamental freedoms	276
bb) Doctrine on emanations of the state	278
d) Summary and discussion	279
aa) General principles	279
bb) Risk of laundered state action in light of the case law analysis	281
(1) Discretionary standards	281
(2) Non-binding demands	283
(3) Meta-regulation	285
(4) Over-blocking risks	288
 <i>III. How can limitations of users' right to freedom of expression be established where public authorities require platform companies to moderate content?</i>	
1. Limitations of users' freedom of expression at the systemic level	291
a) Rational response to platform regulation	292
aa) Coercive nature of platform regulation and rational decision-making	292
bb) Impact of non-regulatory influence modes on regulatory compliance	295
cc) Probabilistic analysis of platform companies' rational response to regulation	300
b) Suitability of speech restrictions as a means to avoid consequences of non-compliance	305
2. Individual limitations of users' freedom of expression	307

<i>IV. To what extent are public authorities obligated to minimize the risks of over-blocking in the context of platform regulation in order to justify limitations of users' right to freedom of expression?</i>	311
1. General principles	312
a) Justified limitations under the Convention	312
b) Justified limitations under the Charter	314
2. Obligation to minimize over-blocking risks	316
a) Case law on safeguards for freedom of expression	317
b) Over-blocking risk factors	320
aa) Vague and abstract substantive provisions	321
bb) Short time frames for content review	325
cc) Monitoring obligations	327
dd) Asymmetric incentives for over-blocking	330
c) Safeguards against over-blocking risks	335
aa) Procedural safeguards	335
(1) Ex ante safeguards	336
(2) Ex post safeguards	339
(a) Transparency	339
(b) Redress	341
(3) Fundamental shortcomings of procedural safeguards	345
bb) Systemic safeguards	347
(1) Formal guarantees and general obligations aimed at protecting users' freedom of expression	348
(2) Freedom of expression by design	349
(3) Fundamental rights impact assessments	353
3. Recommendations on safeguards for freedom of expression	360
<i>V. Conclusion</i>	362
Conclusion	371
Deutsche Kurzzusammenfassung	379
Table of legislation	383
Table of cases	385
Literature	391
Index	429

List of figures and tables

All figures and tables were developed by the author.

List of figures

1. Overview of influences on content moderation	94
2. Impact of the net value of content to platform companies' approach to the limited accuracy of content moderation	98
3. Impact of non-regulatory influences on regulatory compliance	119
4. Scope of DSA	125
5. Effects of platform regulation, systems-level measures, and individual restrictions on users' freedom of expression	212
6. Direct responsibility for platform regulation (causation scenario)	217
7. Vicarious responsibility for content moderation (attribution scenario)	218
8. Factors determining asymmetric incentives (over-blocking risk)	332

List of tables

1. Overview of factors determining the positive or negative value of content	96
2. Overview of content moderation modalities	99
3. Overview of online freedom of expression cases	253
4.1. Negative consequences of (non-)compliance with platform regulation	297
4.2. Negative consequences of (non-)compliance with platform regulation	297
5. Enforcement approaches of different regulatory strategies	301
6. Factors determining the relative costs of non-compliance and of compliance	304

List of abbreviations

AG	Advocate General of the Court of Justice of the European Union
ARIO	ILC Articles on the Responsibility of International Organizations
ARSIWA	ILC Articles on Responsibility of States for Internationally Wrongful Acts
AVMSD	Audiovisual Services Directive
CDA	Communication Decency Act
ChFR	Charter of Fundamental Rights of the European Union
CJEU	Court of Justice of the European Union
DMCA	Digital Millennium Copyright Act
DSA	Digital Services Act Regulation
DSMD	Copyright in the Digital Single Market Directive
EC	European Communities
ECD	E-Commerce Directive
ECHR	European Convention on Human Rights
ECtHR	European Court of Human Rights
EU	European Union
GDPR	General Data Protection Regulation
GIFCT	Global Internet Forum to Counter Terrorism
GNI	Global Network Initiative
ICCPR	International Covenant on Civil and Political Rights
ICJ	International Court of Justice
ILC	International Law Commission
InfoSoc Directive	Directive on the harmonisation of certain aspects of copyright and related rights in the information society.
IP	Internet protocol
KoPI-G	Austrian Communication Platforms Act
NetzDG	German Network Enforcement Act
OECD	Organisation for Economic Co-operation and Development
TERREG	Regulation on addressing the dissemination of terrorist content online
TEU	Treaty on European Union
TFEU	Treaty on the Functioning of the European Union
UN	United Nations
UrhDaG	Act on the Copyright Liability of Online Content Sharing Service Providers

Introduction

I. Setting the stage

Platform companies hold unprecedented power in today's platform society.¹ Just take the deplatforming of Donald Trump by X² and other platform companies in the days following 6 January 2021.³ Given the threat of further violence, this was certainly the right call to make – and arguably long overdue, considering Trump's history of spewing hate and misinformation online.⁴ Still, questions lingered whether this kind of power should be in the hands of private corporations at all, and whether societies should continue to depend on platform companies “doing the right thing if and when they wish, independently of any rules and democratic accountability.”⁵ These concerns only became more pressing once Elon Musk reinstated Donald Trump's X account in November 2022, within weeks of taking over the platform, despite earlier promises not to make “major content decisions or account reinstatements”⁶ before the formation of a new content moderation council.⁷

¹ On the notion of a platform society, see van Dijk, Poell, and de Waal, *The Platform Society*, who refer to “a connective world where platforms have penetrated the heart of societies – disrupting markets and labor relations, transforming social and civic practices, and affecting democratic processes.”; see generally Gillespie, *Custodians of the Internet*, 14, who notes that we are today, by and large, speaking from platforms.

² Previously Twitter.

³ On 6 January 2021, Donald Trump incited a mob to storm the US Capitol in an effort to thwart the election of Joe Biden as US President. For an overview of the actions taken by different companies, see Chrichton, “The Deplatforming of President Trump: A Review of an Unprecedented and Historical Week for the Tech Industry.”

⁴ Cf. Goldberg, “The Scary Power of the Companies That Finally Shut Trump Up”; cf. Kuczerawy, “Does Twitter Trump Trump? A European Perspective.”

⁵ See Floridi, “Trump, Parler, and Regulating the Infosphere as Our Commons,” 1.

⁶ Elon Musk, Tweet, 28 October 2022, <https://twitter.com/elonmusk/status/1586059953311137792> (“Twitter will be forming a content moderation council with widely diverse viewpoints. No major content decisions or account reinstatements will happen before that council convenes.”).

⁷ Cf. Kopps and Katzenbach, “Turning Back Time for Content Moderation? How Musk's Takeover Is Affecting Twitter's Rules.”

Against this background, the normative argument underpinning the thesis is that decisions about what counts as acceptable in the online speech environment should not be taken by platform companies alone. Content moderation, to put it bluntly, is not a private matter.⁸ Instead, public actors with democratic legitimacy need to play a role in setting the objectives of online speech governance. The DSA and other regulatory efforts in the EU are therefore welcome in principle. It is also generally defensible that public authorities are coopting platform companies to restrict illegal online content; coping with the sheer amount of online speech could hardly be accomplished without some degree of public-private cooperation.⁹

Yet, we also need to be mindful of the potential pitfalls of such regulatory interventions. “Determining how and where to regulate [online] speech is among the most important, and most delicate, tasks a government may undertake,” Michael Karanicolas stressed, seeing that “[i]t requires a careful balancing between removing harmful content while providing space for controversial and challenging ideas to be aired, and deterring dangerous speech while minimizing a broader chilling effect that impacts legitimate areas of debate.”¹⁰ Shifting the responsibility to police harmful content to platform companies, as public authorities are increasingly doing, may therefore create its own challenges.¹¹ Although such regulatory approaches are designed to advance public policy objectives, they may actually reduce the authority of state actors in the long term.¹² As David Kaye pointed out, “the pressure on companies has led to an outsourcing of public roles to private actors, which amounts to an expansion of corporate power instead of constraints on it.”¹³

More specifically, delegating the task of online speech control to platform companies carries risks for freedom of expression.¹⁴ As Hannah Bloch-Wehba noted,

“rather than simply compelling intermediaries to delete specific content, governments are foisting upon platforms increasing responsibility for making *legal determinations* regarding speech – a task that might previously have belonged to a court, administrative agency, or other government body accountable to the public.”¹⁵

⁸ Eder, “Making Systemic Risk Assessments Work,” 8.

⁹ Cf. Land, “Against Privatized Censorship,” 395.

¹⁰ Karanicolas, “Squaring the Circle Between Freedom of Expression and Platform Law,” 177 f.

¹¹ Cf. Frosio, “From Intermediary Liability to Responsibility.”

¹² See Land, “Against Privatized Censorship,” 373, fn. 40.

¹³ Kaye, *Speech Police: The Global Struggle to Govern the Internet*, 20.

¹⁴ Cf. Husovec, “(Ir)Responsible Legislature? Speech Risks under the EU’s Rules on Delegated Digital Enforcement,” 3.

¹⁵ Bloch-Wehba, “Global Platform Governance,” 31 f.

In this process, online speech governance risks getting warped into something that does not even remotely resemble traditional speech regulation.¹⁶ Jack Balkin warned that users “get no judicial determination of whether their speech is protected or unprotected when companies block, censor, or take down their speech.”¹⁷ Instead, he continued, “some nontransparent form of private governance or bureaucracy serves as prosecutor, judge, jury, and executioner.”¹⁸ In this context, platform companies have been criticized for their overly vague and opaque content policies and their over-reliance on automation.¹⁹

It has also been extensively argued that platform companies’ motivation to avoid liability, stricter regulation, or administrative fines may incentivize them to err on the side of caution and to systematically restrict legitimate speech.²⁰ Regulatory efforts to address legitimate concerns about unlawful online speech may thus inadvertently erode users’ freedom to impart and receive information and ideas online. In other words, there is a risk of a zero-sum outcome where tackling bad content collaterally hurts good speech.²¹

Moreover, there are concerns – central to this thesis – that delegating speech control responsibilities to platform companies might allow public authorities to bypass their fundamental rights obligations.²² James Boyle recognized this potential more than two decades ago, when he foresaw that states would rely on privatization to regulate the internet:

“[O]ne would want to escape from the practical and legal limitations of a sovereign-citizen relationship. Thus, one might seek out private actors involved in providing Internet services who are not quite as mobile as the flitting and frequently anonymous inhabitants of cyberspace. [...] By enlisting these nimbler, technologically savvy players as one’s private police, one would also gain another advantage: freedom from some of the constitutional and other restraints that would burden the state were it to act directly.”²³

¹⁶ Cf. Balkin, “Free Speech Is a Triangle,” 2028–32, who describes the result of the connection between public and private ordering in online speech governance as “privatized bureaucracy.”

¹⁷ Balkin, 2031.

¹⁸ Balkin, 2031.

¹⁹ See, for instance, Kaye, “Report of the Special Rapporteur on the Promotion and Protection of the Right to Freedom of Opinion and Expression.”

²⁰ Council of Europe, “Comparative Study on Blocking, Filtering and Take-down of Illegal Internet Content,” 30; Keller, “Who Do You Sue? State and Platform Hybrid Power Over Online Speech,” 5; Frosio, “Mapping Online Intermediary Liability,” 26.

²¹ For a critical take on this narrative, see Woods, “Online Harms: Why We Need a Systems-Based Approach Towards Internet Regulation.”

²² See Jørgensen and Pedersen, “Online Service Providers as Human Rights Arbiters,” 180; cf. Land, “Against Privatized Censorship,” 395.

²³ Boyle, “Foucault in Cyberspace: Surveillance, Sovereignty, and Hardwired Censors,” 197.

As Hannah Bloch-Wehba pointed out, this strategy of compelling platform companies to “instantiate and enforce *public* policy preferences” through *private* ordering converts “what might otherwise be private action into heterodox, hybrid public-private governance arrangements in which state and private power are commingled.”²⁴ The resulting blurring of boundaries between private ordering and public regulation makes it difficult to hold public authorities responsible for adverse effects on freedom of expression.²⁵ This raises the question whether public authorities can really break free from “pesky constitutional constraints” when using platform companies as private surrogates.²⁶

II. Hypothesis and research questions

The research hypothesis is that public authorities have an obligation to abstain from violating users’ right to freedom of expression when regulating platform companies’ content moderation processes, regardless of whether or not this regulatory activity directly restricts online speech. In particular, this obligation entails that any limitations of users’ freedom of expression must be justified. The overall research question that motivates this thesis therefore concerns the function of the right to freedom of expression as a constraint on public power: What limits do Art. 10 ECHR and Art. 11 ChFR impose on public authorities’ power to regulate how platform companies moderate online content?

This research question can be broken down into several separate aspects relating, respectively, to content moderation, platform regulation, and platform users’ right to freedom of expression. In developing an answer to this query, several sub-questions, which individually address specific aspects of the overall research question, will guide the analysis. With respect to platform companies’ *content moderation*, the thesis tackles the following questions:

- What does content moderation entail?
- What role does context play in content moderation?
- To what extent does content moderation impact users’ freedom of expression?
- How do different incentives and constraints influence platform companies’ decision-making regarding content moderation?

²⁴ Bloch-Wehba, “Global Platform Governance,” 30.

²⁵ Cf. Jørgensen and Pedersen, “Online Service Providers as Human Rights Arbiters,” 186 f.

²⁶ See Boyle, “A Nondelegation Doctrine for the Digital Age?,” 10.

As regards *platform regulation*, the thesis aims to address the following questions:

- How has platform regulation evolved in the EU?
- What are the rationales behind platform regulation?
- How does platform regulation interact with other influences on platform companies’ decision-making?
- How does platform regulation operate in practice?

Lastly, the thesis explores the following questions regarding users’ *right to freedom of expression*:

- To what extent does platform regulation allow public authorities to launder state action through platform companies’ private ordering?
- To what extent does the case law of the ECtHR and the CJEU on causation and attribution provide avenues for overcoming these risks?
- How can limitations of users’ right to freedom of expression be established where public authorities require platform companies to moderate content?
- To what extent are public authorities obligated to minimize the risks of over-blocking in the context of platform regulation in order to justify such limitations?

III. Research scope

In the following, I will outline the research scope with respect to online platforms, content moderation, platform regulation, and the relevant fundamental rights framework.

1. Content moderation by online platforms

The thesis focusses on content moderation at the application layer of the tech stack. This calls for a more detailed outline of the “what” (content moderation), the “who” (online platforms), and the “where” (application layer).

a) Online platforms

To start with the “who,” the thesis focuses on online platforms for user-generated or user-uploaded content (also referred to as user speech or online speech in the following). Attempts to further define these entities – often vaguely characterized as social media – are made difficult by the ambiguous nature of the term platform, which may also be understood slightly differently in various disci-

plines.²⁷ It follows that “there is no consensus on a single definition of [online platforms], neither in computer science nor in the economic and legal domain.”²⁸ Moreover, regulatory frameworks such as the DSA and the DSMD advanced several new legal concepts related to the notion of online platforms, which reshape and add to pre-existing legal concepts.²⁹

Against this background, the thesis opts for a broad perspective on online platforms. Selectively drawing on existing definitions, online platforms are understood in the thesis as providing three central affordances:

- the technological intermediation of user-generated or -uploaded content,
- the possibility of interactivity among different users and of direct engagement with content,
- the possibility for users to carry out specific activities.³⁰

Naturally, the thesis focusses on online platforms dealing with content moderation questions – although this hardly narrows the scope, since content moderation turns up in unexpected places, from knitting forums³¹ to porn sites³² to the Metaverse.³³ The research scope therefore covers a diverse set of actors, including, but not limited to,

- social media platforms characterized by network connections between users (e. g., Facebook, LinkedIn),
- online media sharing platforms allowing the upload or the livestreaming of content, including images and video (e. g., YouTube, Twitch), music (e. g., Spotify, SoundCloud), and text (e. g., Medium),
- discussion forums allowing users to hold conversations (e. g. Reddit),
- messaging platforms allowing users to communicate and share online content privately (e. g., Telegram),

²⁷ For a concise overview of the term and its history, see Gorwa, “What Is Platform Governance?,” 3; cf. Schwarz, “Platform Logic: An Interdisciplinary Approach to the Platform-Based Economy,” 3; see generally Gillespie, “The Politics of ‘Platforms’” (on the strategic use of the platform metaphor).

²⁸ Bertolini, Episcopo, and Cherciu, “Liability of Online Platforms,” 7 f.

²⁹ Cf. Quintais et al., “Copyright Content Moderation in the EU,” 44–51 (with an instructive overview of the EU law terminology).

³⁰ See DeNardis and Hackl, “Internet Governance by Social Media Platforms,” 2; see also Quintais et al., “Copyright Content Moderation in the EU,” 52 f. (for an overview of different definitions of online platforms proposed in the literature).

³¹ See Copia Institute, “Content Moderation Case Study: Knitting Community Ravelry Bans All Talk Supporting President Trump (2019).”

³² See, for instance, Meineke and Alfering, “We Went Undercover in xHamster’s Unpaid Content Moderation Team.”

³³ See Mak, “I Was a Bouncer in the Metaverse.”

- matchmaking and e-commerce platforms facilitating the transaction of goods and services (e. g., eBay, Airbnb) or the matching between different sets of users (e. g., Tinder),
- platforms for ratings and reviews of third-party services or offerings of different kinds (e. g., Yelp), and
- platforms allowing collaborative production (e. g., Wikipedia).³⁴

For sake of brevity, I will use the term *platform company* to refer to private enterprises offering any kind of platform services and the term *platform* to refer to those services.³⁵ I will not focus on online search engines, unless relevant provisions in the regulatory frameworks analyzed in the thesis specifically refer to them.³⁶

b) Content moderation

As regards the “what,” content moderation (also referred to as content governance, self-regulation, and private ordering³⁷) is “a broad concept with fuzzy borders.”³⁸ It generally refers to platform companies’ practice of setting rules around speech and enforcing them, usually with a mix of human laborers and automated systems at scale.³⁹ In this vein, Kate Klonick characterized content moderation as the “industry term for a platform’s review of user-generated content posted on its site and the corresponding decision to keep it up or take it down.”⁴⁰

The thesis adopts a broader perspective, which is not solely focused on content removal.⁴¹ To this end, it draws on James Grimmelmann’s conception of content moderation; Grimmelmann argued that the term should be understood as the “governance mechanisms that structure participation in a community to facilitate

³⁴ See Bertolini, Episcopo, and Cherciu, “Liability of Online Platforms,” 17; see van Hoboken et al., “Hosting Intermediary Services and Illegal Content Online,” 12–14.

³⁵ See also Gorwa, “What Is Platform Governance?,” 3.

³⁶ See, for instance, Art. 33 DSA.

³⁷ See Gorwa, “The Shifting Definition of Platform Governance.”

³⁸ Quintais et al., “Copyright Content Moderation in the EU,” 33.

³⁹ See Gorwa, “The Shifting Definition of Platform Governance.”

⁴⁰ Klonick, “The Facebook Oversight Board,” 2427.

⁴¹ On this broader perspective, see Douek, “Content Moderation as Systems Thinking,” 531 (“‘Content moderation,’ especially but not exclusively at the largest platforms, now includes many more things than it did even a few years ago: increased reliance on automated moderation; sticking labels on posts; partnerships with fact-checkers; greater platform and government collaboration; adding friction to how users share content; giving users affordances to control their own online experience; looking beyond the content of posts to how users behave online to determine what should be removed; tinkering with the underlying dynamics of the very platforms themselves.”).

cooperation and prevent abuse.”⁴² This definition covers the many ways in which platforms influence users’ online activities,⁴³ including the host of design decisions that structure content flows and user interactions.⁴⁴ It underlines that individuals’ capability for speech, online as offline, is shaped through various and evolving constraints and affordances.⁴⁵ This is echoed by the DSA’s definition of content moderation, which includes

“measures taken that affect the availability, visibility, and accessibility of that illegal content or that information, such as demotion, demonetisation, disabling of access to, or removal thereof, or that affect the ability of the recipients of the service to provide that information, such as the termination or suspension of a recipient’s account.”⁴⁶

Against this background, the term content moderation is used in the thesis to refer to the development of rules regarding online speech, the institutional processes of identifying, adjudicating, and sanctioning content, and, lastly, redress mechanisms that allow users to appeal specific enforcement outcomes.⁴⁷ As in the DSA, this includes algorithmic downranking and other restrictions that stop short of removing content, but only to the extent to which they are used as a means to sanction illegal, prohibited, or unwanted online behavior.⁴⁸ On the other hand, the thesis will not bring demonetization, as another means to sanction content,⁴⁹ into focus. By using the term *content* moderation, I do not mean to exclude restrictions that are not directly applied to individual pieces of content, such as account-level restrictions.

⁴² Grimmelmann, “The Virtues of Moderation,” 47; see also Gorwa, “The Shifting Definition of Platform Governance.”

⁴³ Gorwa, “The Shifting Definition of Platform Governance.”

⁴⁴ Cf. Gorwa; cf. Gorwa, Binns, and Katzenbach, “Algorithmic Content Moderation,” 3.

⁴⁵ Cf. Bietti, “Free Speech Is Circular.”

⁴⁶ Art. 3(t) DSA.

⁴⁷ See, for a similar definition, Gillespie et al., “Expanding the Debate About Content Moderation,” 2.

⁴⁸ De-prioritizing certain pieces of content based on engagement metrics therefore falls outside the scope, whereas demoting content due to its harmful properties is covered by the definition. For a similar discussion on the relation between the concepts of content moderation and content recommending in the DSA context, see Quintais et al., “Copyright Content Moderation in the EU,” 35 f.

⁴⁹ See Llansó et al., “Artificial Intelligence, Content Moderation, and Freedom of Expression,” 18; see generally Caplan and Gillespie, “Tiered Governance and Demonetization: The Shifting Terms of Labor and Compensation in the Platform Economy.”

I will focus on *industrial, commercial* content moderation⁵⁰ by what the DSA refers to as very large online platforms,⁵¹ which typically relies on complex governance structures with tens of thousands of human reviewers as well as large-scale automated systems.⁵² In this sense, content moderation is understood in the thesis “as a project of mass speech administration.”⁵³ This does not only include individual speech restrictions carried out by frontline moderators but also measures at what I refer to as the systems level of content moderation, including design decisions that structure the overall content moderation process.⁵⁴

By contrast, I will not specifically analyze *artisanal* approaches, involving case-by-case review of content by human moderators on a smaller scale (think smaller instances on Mastodon), and *community-reliant* approaches, which commonly combine top-down policy decisions by company staff with a larger group of volunteer human reviewers (think Wikipedia).⁵⁵ I will also primarily focus on content moderation at the individual-firm level. That said, it is important to recognize that there are many overlaps between platform companies’ content moderation processes, from informal exchange between the major US platform companies,⁵⁶ to collective self-regulatory bodies such as the GNI,⁵⁷ to the formation of what Evelyn Douek called content cartels,⁵⁸ such as the joint GIFCT database for terrorist content.⁵⁹

Lastly, three important areas are excluded from the thesis entirely: Practices of data extraction and accumulation associated with the business models of platform companies fall outside the research scope. The thesis also does not comment on platform affordances that allow users to moderate content themselves.⁶⁰

⁵⁰ See Caplan, “Content or Context Moderation?,” 6; see Roberts, “Digital Refuse: Canadian Garbage, Commercial Content Moderation and the Global Circulation of Social Media’s Waste,” 6 f.

⁵¹ On the notion of very large online platforms (and very large online search engines), see Art. 33 DSA.

⁵² See Caplan, “Content or Context Moderation?,” 15–25.

⁵³ On the notion of content moderation bureaucracies, see Douek, “Content Moderation as Systems Thinking.”

⁵⁴ Cf. Douek, 545–48 (on the role of design and affordances in content moderation).

⁵⁵ Caplan, “Content or Context Moderation?,” 15–25.

⁵⁶ See Gillespie et al., “Expanding the Debate About Content Moderation,” 5.

⁵⁷ The GNI was founded in 2008 by a coalition of different platform companies, academics and NGOs with the objective to both coordinate resistance to censorship requests by authoritarian states and to respond to criticisms levied at platforms for accommodating such demands, see Cohen, *Between Truth and Power: The Legal Constructions of Informational Capitalism*, 241.

⁵⁸ Douek, “The Rise of Content Cartels.”

⁵⁹ See below, Ch. 1, Sect. I.3.a.(2)(a).

⁶⁰ On the challenges of allowing political parties to moderate debates in the comments sec-

Lastly, I will not consider online ads as another category of content hosted and moderated by platform companies.

c) Application layer

With respect to the “where,” the thesis focuses on content moderation at the top level of the tech stack,⁶¹ also sometimes referred to as the application layer.⁶² This top level makes up what users experience as the internet, broadly speaking, and it is where platform companies operate.⁶³ Measures taken by actors further down the tech stack – such as cloud service providers, content delivery networks, domain registrars, and internet service providers⁶⁴ – therefore fall outside the research scope.

Still, it is worth noting that there can be an important overlap between the different levels of the tech stack, which may also impact content moderation decisions taken at the top level.⁶⁵ Indeed, we have seen app store providers pressure platform companies and, in some cases, even remove platforms entirely from their ecosystems for a perceived failure to ensure adequate content moderation.⁶⁶ These actions often do not follow clear standards or procedures and tend to be highly opaque,⁶⁷ even though they may have far-reaching consequences.⁶⁸

2. Platform regulation

Platform regulation is used in the thesis to refer to regulation directed at content *moderation* rather than at specific pieces of content. Put differently, the thesis

tions and to hide problematic speech, see generally Kalsnes and Ihlebæk, “Hiding Hate Speech: Political Moderation on Facebook.”

⁶¹ For an overview of the tech stack’s different levels, see Donovan, “Navigating the Tech Stack: When, Where and How Should We Moderate Content?”

⁶² The application layer forms part of the open systems interconnection (OSI) model developed by the International Organization for Standardization, see Cloudflare, What is the OSI model?, <https://www.cloudflare.com/learning/ddos/glossary/open-systems-interconnection-model-osi/>.

⁶³ Cf. Donovan, “Navigating the Tech Stack: When, Where and How Should We Moderate Content?”

⁶⁴ This list is drawn from Joan Donovan’s overview of the different actors at various levels of the tech stack, see Donovan.

⁶⁵ Cf. Gillespie et al., “Expanding the Debate About Content Moderation,” 6 f.

⁶⁶ Cf. Gillespie et al., 7.

⁶⁷ For a first hand account of the influence of Apple and Google on Twitter, see Roth, “I Was the Head of Trust and Safety at Twitter. This Is What Could Become of It.”

⁶⁸ Cf. Kuczerawy, “Freedom of Expression in the Era of Online Gatekeeping,” 284 f. (on the need for safeguards against premature measures by actors down the tech stack).

Index

- Ahmet Yıldırım v. Turkey* 62–63, 229
Automated content moderation 31, 33–45,
51, 93, 327–329
– accuracy 38–40
– bias 40–43
– classifying 36–38
– explainability 43–44
– matching 34–36
- Breyer v. Germany* 227, 286–287
- Cengiz and Others v. Turkey* 72
Chilling effect 79–80, 291
Choice architecture 81, 92, 214
Co-regulation 142–145, 151–153, 334
Code of Conduct on Countering Illegal Hate
Speech Online 145–153
Command-and-control 153–155, 190–195
Community standard *see* content policy
Conditional liability 123–124, 139–141,
332–333
Content moderation
– accuracy 38–40, 91–93, 97–98, 325–326
– appeal mechanism 54, 341–345, 361
– downranking 8, 48, 74–78, 175
– *ex ante* moderation 51–52, 69–70, 138,
327–329, 350
– *ex post* moderation 51–52
– human review 32–33, 44–45
– visibility restriction 46, 48–50, 72–78
Content policy 26–31, 162–163, 322–325
- Data access 52–53, 362
Delfi AS v. Estonia 205, 231–236, 256
Demonetization 8, 50
Digital Services Act *see* DSA
Downgrading *see* downranking
- Downranking 8, 48, 74–78, 175
DSA 111–113, 124–130, 139–141, 170–195
DSMD 130–141, 245–248
Due-diligence obligation 112–113, 180
- Fadeyeva v. Russia* 237–238
Freedom of expression by design 349–353,
361
Fundamental rights impact assessment
353–362
- General monitoring obligation *see* monitor-
ing obligation
Glawischnig-Piesczek 243–245
Governance 20–22
Governing orders 13–14, 371–372
- Human-in-the-loop 344
- Impact assessment *see* fundamental rights
impact assessment
Intermediary liability *see* conditional
liability
Internet referral unit 284
- Kharitonov v. Russia* 230, 317–318, 339
Kotov v. Russia 272–274, 282, 286,
- Liseytseva and Maslov v. Russia* 275
- Meta Oversight Board 54–55
Meta-regulation 16–17, 116–118, 210–213,
252–254, 285–288
Monitoring obligation 242–248, 256–257,
327–329
MTE and Index.hu Zrt v. Hungary 79–80,
233–235, 334–335

- Network Enforcement Act *see* NetzDG
 NetzDG 115–170, 190–192, 323–325, 340
- Out-of-court dispute settlement 177–178, 341–343, 361
- Platform law 27–28
 Platform regulation *see* regulation
Poland v. Parliament and Council 132, 138, 245–248, 319
- Radio France v. France* 268–270
- Regulation
 – co-regulation 142–145, 151–153, 334
 – command-and-control 153–155, 190–195
- conditional liability 123–124, 139–141, 332–333
 – meta-regulation 16–17, 116–118, 210–213, 252–254, 285–288
 Regulatory intermediation 16–17
- SABAM v. Netlog* 69, 232–233, 242–244
SABAM v. Scarlet Extended *see* *SABAM v. Netlog*
 Safe harbor 106–109, 129–130, 139–141
Saliyev v. Russia 271–272, 283
 Shadow banning *see* downranking
 Soft law *see* co-regulation
- Yershova v. Russia* 270–271