

1 **Operational Framework: Institutional Controls - The New Deal** 2 **on Data**

3 Daniel "Dazza" Greenwood^{1,*}, Arkadiusz Stopczynski^{1,2}, Brian Sweatt¹, Thomas Hardjono¹,
4 Alex Sandy Pentland¹

5 **1 MIT**

6 **2 DTU**

7 * **E-mail: dazza@civics.com**

8 **Contents**

9	1 The New Realities of Living in a Big Data Society (Arek)	2
10	2 The New Deal on Data (Arek)	4
11	3 Personal Data: Emergence of a New Asset Class (Thomas)	6
12	4 Enforcing the New Deal on Data (Dazza)	10
13	5 Transitioning End-User Assent Practices (Arek)	13
14	6 Business, Legal, and Technical Dimensions of Big Data Systems (Dazza)	14
15	7 Big Data and Personal Data Institutional Controls (Thomas)	16
16	8 Scenarios of Use in Context (Dazza)	20
17	8.1 Example Scenario: Research Systems	21
18	8.2 Scenarios of Use Today, Tomorrow and the Day After	23
19	9 Future Research (Brian)	25
20	9.1 Research on Design and Deployment of Big Data Systems	26
21	9.2 Research on Big Data for Design of Institutions	28

23 1 The New Realities of Living in a Big Data Society (Arek)

24 To realize the promise and prospects of a Big Data society and avoid its security and confiden-
25 tiality perils, institutions are updating operational frameworks governing business, legal, and
26 technical dimensions of their internal organization and interactions with the outside world. In
27 this chapter we explore the emergence of the Big Data society, outline ways to support it in the
28 context of institutional controls within the framework of New Deal on Data, and describe future
29 directions of research and development.

30 The control points traditionally relied upon as part of corporate governance, management
31 oversight, legal compliance, and enterprise architecture must evolve and expand to match oper-
32 ational frameworks for Big Data. An operational framework used for a Big Data-driven organi-
33 zation requires a balanced set of institutional controls. These must support and reflect greater
34 user control over personal data, and large scale interoperability for data sharing between and
35 among institutions. Core capabilities of these controls include responsive rule-based systems
36 governance and fine-grained authorizations for distributed rights management.

37 Sustaining a healthy, safe, and efficient society is a scientific and engineering challenge go-
38 ing back to the 1800s, when the Industrial Revolution spurred rapid urban growth, creating
39 huge social and environmental problems. The remedy then was to build centralized networks
40 that delivered clean water and safe food, enabled commerce, removed waste, provided energy,
41 facilitated transportation, and offered access to centralized healthcare, police, and educational
42 services. Those networks formed the backbone of the society as we know it today.

43 These century-old solutions are however becoming increasingly obsolete and inefficient. We
44 have cities jammed with traffic, world-wide outbreaks of disease that are seemingly unstoppable,
45 and political institutions that are deadlocked and unable to act. We face the challenges of global
46 warming, uncertain energy, water, and food supplies, and a rising population and urbanization,
47 that will add 350 million people to the urban population by 2025 in China alone [14].

48 It does not have to be this way. We can have cities that are protected from pandemics, energy
 49 efficient, have secure food and water supplies, and have much better government. To reach these
 50 goals, however, we need to radically rethink our approach. Rather than static fixed systems,
 51 separated by function — water, food, waste, transport, education, energy — we must consider
 52 them as dynamic, data-driven networks. Instead of focusing only on access and distribution,
 53 we need the networked and self-regulating systems, driven by the needs and preferences of the
 54 citizens.

55 Sustainable future society depends on using our new technologies to create a *nervous system*
 56 maintaining the stability of government, energy, and public health systems around the globe.
 57 The digital feedback technologies are today capable of creating a level of dynamic responsiveness
 58 that our larger, more complicated modern society requires. We must reinvent the systems of the
 59 societies within a control framework: sensing the situation, combining these observations with
 60 models of demand and dynamic reaction, and finally using the resulting predictions to tune the
 61 system to match the demands.

62 The engine driving this nervous system is Big Data – the newly ubiquitous digital data,
 63 now available about all aspects of human life. We can analyze patterns of human experience
 64 and ideas exchange within the *digital breadcrumbs* that we all leave behind as we move through
 65 the world: call records, credit card transactions, GPS location fixes, among others [24]. By
 66 recording our choices, these data tell the story of our lives. And this may be very different from
 67 what we decide to put on Facebook or Twitter; our postings there are what we choose to tell
 68 people, edited according to the standards of the day and filtered to match the persona we are
 69 building. Mining social networks can give some great insights about human nature [4, 28, 42];
 70 who we really are is however even more accurately determined by where we spend our time and
 71 which things we buy, rather than just what we say we do [27].

72 The process of analyzing the patterns within these digital breadcrumbs is called reality
 73 mining [13, 32], and through it we can learn an enormous amount about who we are. The
 74 Human Dynamics research group at MIT have found that we can use them to tell if we are

likely to get diabetes [33], or whether we are the sort of person who will pay back loans [34]. By analyzing these patterns across many people, we are discovering that we can begin to explain many things — crashes, revolutions, bubbles — that previously appeared to be random acts of God [30]. For this reason the magazine *Technology Review* named our development of reality mining as one of the ten technologies that will change the world [17].

2 The New Deal on Data (Arek)

The digital breadcrumbs we leave behind provide clues about who we are, what we do and want. This makes these personal data — data about individuals — immensely valuable, both for public good and for private companies. As European Consumer Commissioner, Meglena Kuneva said recently, “Personal data is the new oil of the Internet and the new currency of the digital world” [23]. This new ability to see the details of every interaction can be however used for good or for ill. Therefore, maintaining protection of personal privacy and freedom is critical to our future success as a society. We need to enable even more data sharing for the public good; at the same time, we need to do a much better job in protecting the privacy of the individuals.

A successful data-driven society must be able to guarantee that our data will not be abused; perhaps especially that government will not abuse the power conferred by access to such fine-grain data. The abuses may be directly targeted at users, for example by offering them higher insurance rates based on their shopping history [16], or create problems for the entire society in longer run, for example by limiting user choices and closing them into information bubbles [19]. To achieve the positive possibilities of the new society, we require the *New Deal on Data*, workable guarantees that the data needed for public good are readily available while at the same time protecting the citizenry [32].

The key insight that motivates the idea of the New Deal on Data is that our data are worth more when shared, because these aggregated data — averaged, combined across population, and often distilled to high-level features — inform improvements in systems such as public health, transportation, and government. For instance, we have demonstrated that data about the way

101 we behave and where we go can be used to minimize the spread of infectious disease [26,33]. Our
102 research has reported how we were able to use these digital breadcrumbs to track the spread of
103 influenza from person to person on an individual level. And if we can see it, we can stop it.

104 Similarly, if we are worried about global warming, these shared, aggregated data can show
105 us how patterns of mobility relate to productivity [31]. It can provide us with the ability to
106 design cities that are more productive and, at the same time, more energy efficient. But in order
107 to obtain these results and make a greener world, we need to be able to see the people moving
108 around; this depends on many people willing to contribute their data, even if only anonymously
109 and in aggregate.

110 To enable sharing of personal data and experiences, we need secure technology and regulation
111 that allow individuals to safely and conveniently share personal information with each other,
112 with corporations, and with government. Consequently, the heart of the New Deal on Data
113 must be to provide both regulatory standards and financial incentives that entice owners to
114 share data, while at the same time serving the interests of both individuals and society at large.
115 We must promote greater idea flow among individuals, not just corporations or government
116 departments.

117 Unfortunately, today most personal data are siloed off in private companies and therefore
118 largely unavailable. Private organizations collect the vast majority of the personal data in the
119 form of mobility patterns, financial transactions, phone and Internet communications. These
120 data must not remain the exclusive domain of private companies, because then they are less
121 likely to contribute to the common good. These private organizations must be thus the key
122 players in the New Deal on Data framework for privacy and data control. Likewise, these data
123 should not become the exclusive domain of the government, as this will not serve the public
124 interest of transparency; we should be suspicious of trusting the government with such power.
125 Ultimately, the entities who should be empowered to share and make decisions about their data,
126 are people themselves: users, participants, citizens.

127 Through the years, the great goal of human societies was to find the efficient ways of gov-

ernance. The Big Data transformation can contribute to this ultimate goal of providing the society with tools to analyze and understand what needs to be done, and to reach the consensus on how to do it. This goes beyond just creating more communication platforms; the assumption that more interactions between users will result in better decisions being made, may be very misleading. Although in the recent years we have seen some great examples of using social networks for better organization in society, for example during political protests [6, 18], we are not even close to the point where we can start reaching consensus about the big problems: epidemics, climate change, pollution. We can improve the discussions by making them data driven, involving both experts and wisdom of the crowds – users themselves interested in improving the society. The problems we are dealing with as a now global society are more difficult than ever. We are responsible for many of them, and being able to tackle them on a global scale is necessary for our, mankind, survival.

3 Personal Data: Emergence of a New Asset Class (Thomas)

It has long been recognized that the first step to promoting liquidity in land and commodity markets is to guarantee ownership rights so that people can safely buy and sell. Similarly, the first step toward creating greater idea and idea flow (‘idea liquidity’) is to define ownership rights. The only politically viable course is to give individual citizens rights over data that are about them and in fact, in the European Union these rights flow directly from the constitution **AS: Citation? There is no ‘EU constitution’ per se.** . We need to recognize personal data as a valuable asset of the individual that is given to companies and government in return for services.

The simplest approach to defining what it means to own your own data is to draw an analogy with the English common law ownership rights of possession, use, and disposal:

- You have the right to possess data about you. Regardless of what entity collects the data, the data belong to you, and you can access your data at any time. Data collectors thus

153 play a role akin to a bank, managing the data on behalf of their customers.

154 • You have the right to full control over the use of your data. The terms of use must be opt-
 155 in and clearly explained in plain language. If you are not happy with the way a company
 156 uses your data, you can remove the data, just as you would close your account with a bank
 157 that is not providing satisfactory service.

158 • You have the right to dispose of or distribute your data. You have the option to have data
 159 about you destroyed or redeployed elsewhere.

160 Individual rights to personal data must be balanced with the need of corporations and govern-
 161 ments to use certain data-account activity, billing information, and so on-to run their day-to-day
 162 operations. This New Deal on Data therefore gives individuals the right to possess, control, and
 163 dispose of copies of these required operational data, along with copies of the incidental data
 164 collected about you such as location and similar context.

165 Note that these ownership rights are not exactly the same as literal ownership under modern
 166 law, but the practical effect is that disputes are resolved in a different, simpler manner than
 167 would be the case for (as an example) land ownership disputes.

168 In 2007, one author (Pentland) first proposed the New Deal on Data to the World Economic
 169 Forum [43]. Since then, this idea has run through various discussions and eventually helped
 170 shape the 2012 Consumer Data Bill of Rights in the United States, along with a matching
 171 declaration on Personal Data Rights in the EU. These new regulations hope to accomplish the
 172 combined trick of breaking data out of the current silos, thus enabling public goods, while at
 173 the same time giving individuals greater control over data about them. But, of course this is
 174 still a work in progress and the battle for individual control of personal data rages onward.

175 The World Economic Forum (WEF) has dubbed personal data as the “New Oil” or resource
 176 of the 21st century [43]. The discovery of oil and the subsequent development of the oil industry
 177 over the past 100 years has spurred not only the development of the automobile industry but also
 178 the creation of the global transportation infrastructure, including the massive freeway networks

179 that we see today in the developed nations. The “personal data sector” of the economy today is
 180 still in its infancy, its state akin to the oil industry at the late 1890s prior to the development of
 181 the Model-T Ford automobile. The productive collaboration between the Government (building
 182 the state owned freeways), the private sector (mining and refining oil, building automobiles) and
 183 the citizen (the user-base of these services) allowed the develop nations to expand its economies
 184 by creating new markets adjacent to the automobile and oil industries.

185 If personal data, as the new oil, is to reach its global economic potential, there needs to be
 186 a productive collaboration between all the stakeholders in the establishment of a *personal data*
 187 *ecosystem*. As mentioned in [43], a number of fundamental questions about privacy, property,
 188 global governance, human rights — essentially around who should benefit from the products
 189 and services built upon personal data — are major uncertainties shaping the opportunity. The
 190 rapid rate of technological change and commercialization in using personal data is undermining
 191 end user confidence and trust.

192 The current personal data ecosystem is fragmented and inefficient. Too much leverage is
 193 currently being accorded to service providers that on-board and register end-users. These siloed
 194 repositories of personal data exemplifies the fragmentation of the ecosystem. These repositories
 195 contain data of varying qualities. Some are attributes of persons that are unverified, while
 196 other represent higher quality data that have been cross-correlated with other data points of the
 197 end-user.

198 For many participants, the risks and liabilities exceed the economic returns. Besides not
 199 having the infrastructure and tools to manage personal data, many end-users simply do not see
 200 the benefit of fully participating in the ecosystem. The current focus of many Internet-based
 201 service providers is to capture as much personal data from the end-user and to sell this data into
 202 the advertising industry. Personal privacy concerns are thus inadequately addressed at best,
 203 or simply overlook in the majority of the cases. The current technologies and laws fall short
 204 of providing the legal and technical infrastructure needed to support a well-functioning digital
 205 economy.

206 Recently, we have shown how challenging, but also feasible, it is to open such institutional
 207 Big Data. In the Data For Development (D4D) Challenge <http://www.d4d.orange.com/home>,
 208 the telecom operator Orange opened access to a large dataset of call detail records (CDRs) from
 209 the Ivory Coast. Working with the data as part of a challenge, teams of researchers came up
 210 with life-changing insights for the country. For example, one team developed a model for how
 211 disease spread in the country and demonstrated that information campaigns based on one-to-one
 212 phone conversations among members of social groups can be an effective countermeasure [25]. In
 213 releasing and analysing this data, the privacy of the people who generated the data was protected
 214 not only by the technical means, such as removal of the Personally Identifiable Information
 215 (PIIs), but also by legal means, with the researchers signing an agreement they will not use the
 216 data for re-identification or other nefarious purposes. As we have seen in several cases, such as
 217 the Netflix Prize privacy disaster [29] and other similar privacy breaches [37], true anonymization
 218 is extremely hard. In the Unique in the Crowd [10], de Montjoye et al. showed that even though
 219 human beings are highly predictable [35], we are also very unique. Having access to one dataset,
 220 it may be easy to uniquely fingerprint someone based on just few datapoints, and use this
 221 fingerprint to discover their true identity. The higher the resolution of the data, the easier it
 222 gets to identify a person from this type of data.

223 The report of the World Economic Forum [43] also suggest a way forward by recommending
 224 a number of areas where efforts could be directed:

- 225 • Alignment of key stakeholders: Citizens, the private sector and the public sector need to
 226 work in support of one another. Efforts such as NSTIC [38] — albeit still in its infancy —
 227 represents a promising direction for a global collaboration.
- 228 • Viewing “data as money”: There needs to be a new change in mindset where an individual’s
 229 personal data items are viewed and treated in the same way as their money. These personal
 230 data items would reside in an “account” (like a bank account) where it would be controlled,
 231 managed, exchanged and accounted for just like personal banking services operate today.

- End-user centricity: All entities in the ecosystem need to recognize that end-users are vital and independent stakeholders in the co-creation and value exchange of services and experiences. Efforts such as the *User managed Access* (UMA) initiative [2] point in the right direction by designing systems that are user-centric and managed by the user.

Opening data from the silos by publishing static datasets — collected at some point and unchanging — is important, but it is only the first step. We can do even more substantial things when the data is available in real time and can become part of a society’s nervous system. Epidemics can be monitored and prevented in real time [33], underperforming students can be helped, and people with health risks can be treated before they get sick [9]. The same data can potentially be used for stalking, burglarizing one’s home, and as justification to charge people more for an insurance policy.

4 Enforcing the New Deal on Data (Dazza)

How can we enforce this New Deal? The threat of legal action alone is important, but insufficient, because if you cannot see abuses then you cannot prosecute them. Moreover, who wants more lawsuits anyway? Enforcement can be addressed in significant ways without prosecution of public statute or regulation at all. In many fields, companies and governments rely upon multi-party frameworks of agreed rules governing common business, legal, and technical practices to create effective self-organization and enforcement. These approaches hold promise as a method for using institutional controls to form a reliable operational framework balancing the needs for big data, privacy, and access.

One current best practice is a system of data sharing called trust networks. Trust networks are a combination of networked computers and legal rules defining and governing expectations regarding data. With respect to data belonging to individuals, these networks of technical and legal rules keeps track of user permissions for each piece of personal data, and a legal contract that specifies both what you can and cannot do with the data and what happens if there is a

violation of the permissions. For example, in such a system all personal data can have attached labels specifying what the data can and cannot be used for. These labels are exactly matched by the network's system rules and terms in legal contracts between all the participants, stating penalties for not obeying the permission labels. These rules can, and often do, reference or require audits of relevant systems and data use, demonstrating how traditional internal controls can be leveraged as part of the transition to more novel trust models.

Complete tracking and regulation of every aspect of a trust network is not the goal or even desirable in order to achieve effective enforcement. Rather, the rules for a trust network align enforcement with the highest priority issues and those upon which trust of participants is premised. The relevant issues arise from the dynamics of data flows, underlying trust models, and contextual scenarios within which the networked data and the relationships of parties in the trust network **AS: This sentence is hard to understand. Missing verb?** . When a trust network involves use of personal data, then the user permissions and corresponding limits on use are fundamental to the trust model. In this context, the permissions, including the provenance of the data, should require appropriate levels of audit. A well designed trust network, elegantly integrating computer and legal rules, allows automatic auditing of data use and allows individuals to change their permissions and withdraw data.

Having system rules applicable to the networks, applications, and data as well as all the services providers other intermediaries, and the users themselves is the mechanism for establishing and operating a trust network. System rules are sometimes called operating regulations in the credit card context, or known as trust frameworks in the identity federations context, or trading partner agreements in a supply value chain context. There are many general examples of multiparty shared architectural and contractual rules that share the generic characteristic of creating binding obligations and enforceable expectations on all participants in scalable networks. Another common characteristic of the system rules design pattern is that the participants in the network can be widely distributed across very heterogeneous business ownership boundaries, legal governance structures, and technical security domains. Yet, the parties need not agree to

conform all or most aspects of their basic roles, relationships, and activities in order to connect to to systems of a trust network. Cross-domain trusted systems must, by their nature, focus mandatory and enforceable rules narrowly upon the critical items that must be commonly agreed in order for that network to achieve it's purpose.

For example, institutions participating in credit card and automated clearinghouse debit transactional networks are subject to profoundly different sets of regulations, business practices, economic conditions, and social expectations. The network rules focus upon the topmost agreed items affecting interoperability, reciprocity, risk, and revenue allocation. The knowledge that fundamental rules are subject to enforcement actions is one of the foundations of trust as well as a motivation to prevent or address violations before they trigger penalties. A clear example of this approach can be found with the Visa Operating Rules, covering a vast global real-time network of parties that agree to rules governing their roles in the system as merchants, banks, transaction processors, individual or business card holders, and other key system roles.

A system like this has made the interbank money transfer system among the safest systems in the world and the daily backbone for exchanges of trillions of dollars, but until recently such systems were only for the 'big guys'. To give individuals a similarly safe method of managing personal data, the Human Dynamics research group at MIT, in partnership with the Institute for Data Driven Design, co-founded by John Clippinger and one author (Pentland), have helped build open Personal Data Store (openPDS) [11]. See <http://openPDS.media.mit.edu> for project information and <https://github.com/HumanDynamics/openPDS> for the open source code.

The openPDS is a consumer version of a personal cloud trust network that we are now testing with a variety of industry and government partners. Soon, sharing your personal data could become as safe and secure as transferring money between banks.

The Human Dynamics Lab has applied the system rules approach to development of integrated business, technical architecture, and rules large scale institutional use of personal data stores, available as an example under MIT's creative commons license by MIT, at <https://github.com/HumanDynamics/openPDS>

311 `//github.com/HumanDynamics/SystemRules.`

312 The capacity to apply the appropriate methods of enforcement for a trust network depend
 313 upon a clear understanding and agreement among parties about the purpose of the trusted
 314 system and the respective roles or expectations of those connecting as participants. Therefor,
 315 an anchor is needed to a clear context of a Big Data operational framework and institutional
 316 controls appropriate for access and confidentiality or privacy. The following section posits the
 317 trust model and signature traits of such a context, through the lens of the New Deal on Data.

318 **5 Transitioning End-User Assent Practices (Arek)**

319 The way users grant authorizations to their data is not a trivial matter. The flow of personal
 320 information, such as location data, purchases, health records can be very complex. Every tweet,
 321 geo-tagged picture, phone call, or purchase with credit card, provide the user's location not only
 322 to the primary service, but also to all the applications and services that have been authorized
 323 to access and re-use these data. The authorizations may come from the end-user or be granted
 324 by the collecting service, based on an umbrella terms of service, allowing the re-use of the data.
 325 Implementation of such flows was a crucial part of the Web 2.0 revolution, realized with RESTful
 326 APIs, mashups, and authorization-based access. The way the personal data travel between the
 327 services has however become arguably too complex for a user to handle and manage.

328 Increasing the amount of data controlled by the user and granularity of this control is mean-
 329 ingless if it cannot be exercised in an informed way. For many years, the End User License
 330 Agreements (EULAs), long incomprehensible texts have been accepted blindly by the user,
 331 trusting they have not agreed to anything that could harm them. The process of granting the
 332 authorizations cannot be too complex, as it would prevent the user from understanding her deci-
 333 sions. At the same time, it cannot be too simplistic, as it may not sufficiently convey the weight
 334 of the privacy-related decisions. It is a challenge in itself, to build the end-user assent systems
 335 that allow the user to understand and adjust their privacy settings. Complex EULAs do not
 336 promote the privacy of the users, effectively pushing them to press *I Agree* in every presented

337 window.

338 This gap between the interface — single click — and the effect, can render the data owner-
339 ship meaningless; the click may wrench people and their data into systems and rules that are
340 antithetical to fair information practices, such as is prevalent with today's end-user licenses in
341 cloud services or applications. Managing the potentially long term and opposite dynamics fueled
342 by old deal systems operating simultaneously with the new deal systems is an important design
343 and migration challenge during the transition to a Big Data economy. During this transition
344 and after the New Deal on Data is no longer new, personal data must continue to flow in order
345 to be useful. Protecting the data of people outside of the user-controlled domain is very hard
346 without a combination of cost effective and useful business practices, legal rules, and technical
347 solutions.

348 We envision Living Informed Consent, where the user is entitled to know what data is being
349 collected about her by which entities, empowered to understand the implications of data sharing,
350 and finally put in charge of the sharing authorizations. We suggest the readers ask themselves a
351 question: *Which services know which city I am in today?*. Google? Apple? Twitter? Amazon?
352 Facebook? Flickr? This small application we have authorized a few years ago to access our
353 Facebook check-ins and forgot since then? This is an example of a fundamental question related
354 to user privacy and assent, and yet finding the answer to it may be surprisingly difficult in today's
355 ecosystem. We can hope that most of the services treat the data responsibly and according to
356 user authorizations. In the complex network of data flows however, it is relatively easy for the
357 data to leak to services careless with it or simply malicious [7]. We need to build the solutions
358 to help the user to make well thought-through decisions about data sharing.

359 6 Business, Legal, and Technical Dimensions of Big Data Sys- 360 tems (Dazza)

361 When it comes to data intended to be accessible over networks — whether big, personal, or
362 otherwise — the traditional container of an institution makes less and less sense. Institutional
363 controls apply, by definition by or to some type of institutional entity such as a business, gov-
364 ernmental, or religious organization. A combined view of the business, legal, and technical facts
365 and circumstances surrounding big data is necessary to know what access, confidentiality, and
366 other expectations exist. The relevant contextual aspects of Big Data of one institutional is often
367 profoundly different from that of another. As more and more organizations use and rely upon
368 big data, a single formula for institutional controls will not work for increasingly heterogeneous
369 business, legal and technical environments in play.

370 Looking at an institution as a business, legal, and technical ‘system’ is one effective approach
371 for dealing with the inherent complexity of managing heterogeneous and distributed networks of
372 actors and interactions. The business models, interface-point operational practices and relevant
373 assumptions must be consistent and frequently carefully agreed upon at an executive level by
374 and with institutions as part of the value exchange involving data and access to high value,
375 mission critical or sensitive systems and services. The applicable legal frameworks, common
376 assumptions regarding likely allocation of liability and resolution of disputes in the event of
377 losses, and expected types of contracting practices need to reflect and support the business
378 goals and purposes for the system and data. When technical standards are selected, configured
379 and applied to systems they too must support and reflect the business and legal dimensions and
380 be supported and reflected by those dimensions.

381 Once a systems view is adopted, there is a tractable starting point to narrow or broaden
382 the scope of view to see the smaller and larger systems and to make better and more effective
383 use and control of big data. Within a given institution, there may in fact be many different
384 discernable institutions and corresponding systems and any given system of one institution will

385 frequently in fact exist across many different discernable institutions. However, defining as a
 386 ‘system’ the thing to which institutional controls apply provides an achievable and measurable
 387 basis for balancing privacy, access and other interests in big data. **AS: The paragraph above**
 388 **is hard to understand I think.**

389 Many organizations are structured with clear leadership on business, legal, and technical
 390 issues functionally assigned to top level executive roles. Business issues are typically allocated
 391 to roles such as CEO, COO or CFO, while leadership on legal issues is commonly assigned to
 392 roles like general counsel and regulatory compliance and technical leads are often the roles of
 393 CIO, CTO or CSO. Having top level leadership for each of the business, legal, and technical
 394 aspects of a trust network is a critical success factor.

395 7 Big Data and Personal Data Institutional Controls (Thomas)

396 The phrase “institutional controls” refers to safeguards and protections by use of legal, policy,
 397 governance, and other non-strictly technical, engineering, or mechanical measures. The phrase
 398 institutional controls in a Big Data context can perhaps best be understood by examining how
 399 the concept has been applied to other domains. The most prevalent use of institutional controls
 400 has been in the field of environmental regulatory frameworks.

401 A good example of how this concept supports and reflects the goals and objectives of en-
 402 vironmental regulation can be found in the policy documents of the Environmental Protection
 403 Agency (EPA). This following definition is instructive, and is part of the Institutional Control
 404 Glossary of Terms [40]:

405 “Institutional Controls - Non-engineering measures intended to affect human activi-
 406 ties in such a way as to prevent or reduce exposure to hazardous substances. They
 407 are almost always used in conjunction with, or as a supplement to, other measures
 408 such as waste treatment or containment. There are four categories of institutional
 409 controls: governmental controls; proprietary controls; enforcement tools; and infor-

410 mational devices.”

411 Going deeper, the article by DeMeo and Doar [12] defines institutional controls thusly:

412 “Institutional controls are administrative and legal controls that help minimize the
413 potential for human exposure to contamination and/or protect the integrity of the
414 physical remedy. They can include recorded restrictive covenants, but land use
415 laws and regulations, deed restrictions, department consent orders, and conservation
416 easements are all institutional controls.”

417 In domains of information technology, this approach is most commonly reflected as “enter-
418 prise controls” related to security. See, for example, the report [22] stating: “Enterprise mobility
419 technologies, especially those designed to retrofit enterprise controls on top of consumer mobile
420 devices, are rapidly evolving. This was a message we heard loud and clear in the study.” This
421 study and analysis also reveals much about the internal controls needed to accommodate mobile
422 device use by employees. In both capacities as employee, consumer, and other roles, the use of
423 mobile devices triggers myriad legal, policy, and other implications for institutional controls.

424 In the legal domain, this concept frequently emerges under the moniker “regulatory compli-
425 ance” or “legal compliance” anchored in legal and regulatory frameworks such as Health Insur-
426 ance Portability and Accountability Act (HIPAA) and Sarbanes-Oxley (SOX). These statutory
427 legal frameworks require covered organizations to established integrated sets of governance,
428 legal, transactional, security, and other internal controls to avoid violating the rules. The in-
429 stitutional controls are accomplished in tight integration with engineering and other measures
430 in order to ensure compliance and to control legal and security risk. The use of institutional
431 controls of this type are fundamental methods for achieving and maintaining the transition to a
432 digital, networked, and Big Data footing for any private company, government agency, or other
433 organization.

434 Consider again the analogy of institutional controls in the context of environmental law, and
435 how these types of measures can be applied in the Big Data, privacy, and access context to digital

environments. Given the relatively mature and stable state of environmental regulation, there is much to be learned by examining this context of institutional controls. Environmental regulatory compliance with waste management cleanup requirements could include institutional controls restricting land use on adjacent property. In these situations, it is possible that the remediation strategy requires significant use of land outside the property boundaries of the cleanup site. In these cases, the regulators and the land owner responsible for the regulated property must find ways to ensure a common approach among multiple owners and across multiple property environments. Use of measures such as a clauses on the relevant deeds, an enforceable consent order, or regulations and zoning rules are examples of more severe institutional controls that can be employed to ensure consistent and effective actions are taken across ownership and real property boundaries.

See, for example, Florida Department of Environmental Protection (FDEP), Division of Waste Management [15] which states that “...RMO III does contemplate contamination beyond the Property boundaries, which would require agreement by the adjacent owners to put an RC on their properties as well.”

The concept of an “institutional control boundary” is especially clarifying and powerful when applied to the networked and digital boundaries of an institution. In the context of Florida’s environmental regulation frameworks, the phrase is applied to describe the various types of combinations risk management levels related to target cleanup standards and extend beyond the area of a physical property boundary. Also see a recent University of Florida report on Development of Cleanup Target Levels (CTLs) [8] stating “Risk Management Options Level III, like Level II, allows concentrations above the default groundwater CTLs to remain on site. However, in some rare situations, the institutional control boundary at which default CTLs must be met can extend beyond the site property boundary.”

The EPA provides considerable information on the nature and use of institutional controls, including situations when the situational scope extends to adjacent properties owned by third parties. See, generally, *EPA Hazardous Waste Corrective Action Guidance on Institutional Con-*

463 trols [40]. Also see: *Institutional Controls Bibliography: Institutional Control, Remedy Selection,*
 464 *and Post-Construction Completion Guidance and Policy, December 2005* [39].

465 When institutional controls would apply to “separately owned neighboring properties” a
 466 number of issues arise. Engagement with affected third parties, requiring the party responsible
 467 for site cleanup to use “best efforts” to attain agreement by third parties to institute the relevant
 468 institutional controls, use of third party neutrals to resolve disagreements regarding the appli-
 469 cation with institutional control,s or forcing an acquisition of the neighboring land by forcing
 470 the party responsible to purchase the property of by purchase of the property directly by the
 471 EPA [41].

472 In the context of Big Data, privacy, and access, institutional controls are seldom, if ever,
 473 the result of government regulatory frameworks such as are seen in the environmental waste
 474 management oversight by the EPA. Rather, institutions applying measures constituting institu-
 475 tional controls in the big data and related information technology and enterprise architecture
 476 contexts will typically employ governance safeguards, business practices, legal contracts, techni-
 477 cal security, reporting, and audit programs and a various risk management measures. Inevitably,
 478 institutional controls for Big Data will have to operate effectively across institutional boundaries,
 479 just as environmental waste management internal controls must sometimes be applied across real
 480 property boundaries and may subject multiple different owners to enforcement actions corre-
 481 sponding to the applicable controls. Short of government regulation, the use of system rules as a
 482 general model are one widely understood, accepted, and efficient method for defining, agreeing,
 483 and enforcing institutional and other controls across business, legal, and technical domains of
 484 ownership, governance, and operation.

485 The use of system rules and integrated participation agreements by developers and end-
 486 users is a way to ensure intended operational frameworks conform to applicable institutional
 487 controls. The example of Living Informed Consent described in this chapter, demonstrates how
 488 institutional controls comprised of legal and definite workflow measures, in concert with technical
 489 methods, can result in a higher level of performance, while appropriately balancing legitimate

490 interests of various parties regarding use and access to personal data.

491 Following the World Economic Forum recommendations of treating personal data stores in
 492 the manner of bank accounts [43], there are a number of infrastructure improvements that need to
 493 be realized, if the personal data ecosystem is to flourish and deliver new economic opportunities.
 494 We believe the following infrastructure improvements are necessary for the coming personal data
 495 ecosystem: **AS: We should remove the bullets, turn them into continuous text.**

- 496 • *New global data provenance network*: In order for personal data to be treated like bank
 497 accounts, the origin information regarding data items coming into the data store must be
 498 maintained [21]. In other words, the provenance of all data items must be accounted for
 499 by the IT infrastructure upon which the personal data store operates. The heterogeneous
 500 provenance databases must then be interconnected in order to provide a resilient and
 501 scalable platform for audit and accounting systems to track and reconcile the movement
 502 of personal data from the respective data stores.
- 503 • *Trust network for computational law*: In order for trust to be established between parties
 504 who wish to exchange personal data, we foresee that some degree of “computational law”
 505 technologies may have to be integrated into the design of personal data systems. Such
 506 technologies should not only verify terms of contracts (e.g. terms of data use) against user-
 507 defined policies but also have mechanisms built-in to ensure non-repudiation of entities
 508 who have accepted these digital contracts. Efforts such as [1, 2] are beginning to bring
 509 non-repudiation and enforceability of contracts into the technical protocol flows.
- 510 • *Development of institutional controls for digital institutions*: Currently there are a number
 511 of proposal for the creation of virtual currencies (e.g. BitCoin [5], Ven [36]) in which the
 512 systems have the potential to evolve into self-governing “digital institutions” [20]. Such
 513 systems and institutions that operate on them will necessitate the development of a new
 514 paradigm to understand the aspects of institutional control within their context.

515 8 Scenarios of Use in Context (Dazza)

516 Supporting the effective development of institutional controls for big data requires an under-
 517 standing of how to define and work with the applicable context surrounding the scenarios within
 518 which the Big Data exists. In particular, the New Deal on Data will require a set of Institu-
 519 tional Controls involving governance, business, legal, and technical aspects that are knowable
 520 only with reference to the relevant context of a factually based scenario of use. The following
 521 scenarios demonstrate signature features of the New Deal on Data in various contexts and serve
 522 as an anchor to evaluate what Institutional Controls are well aligned.

523 8.1 Example Scenario: Research Systems

524 **AS: This entire section requires significant write-through.**

525 Computational Social Science (CSS) studies are based on data collected often with an ex-
 526 tremely high resolution and scale [24]. Using computational power combined with mathematical
 527 models, such data can be used to provide insights into human nature. Much of the data collected,
 528 for example mobility traces are sensitive and private; most individuals would feel uncomfortable
 529 sharing them publicly. The need for solutions to ensure the privacy of the individuals has grown
 530 alongside the data collection efforts.

531 The data collection in the CSS context is based on the informed consent of the partici-
 532 pants. Countries have different bodies regulating such studies, for example Institutional Research
 533 Boards (IRBs) in the US. Although certain minimal requirements for implementing informed
 534 consent exist**AS: reference**, they are often not very well suited for the large-scale studies,
 535 where the amount and sensitivity of the data calls for sophisticated privacy controls. As the
 536 scale of the studies grows, in terms of the number of participants, collected bits per user, and
 537 duration, the EULA-style informed consent is no longer sufficient and makes it hard to claim
 538 that participants in fact expressed informed consent.

539 One author (Stopczynski) deployed this year a 1,000 phones study at Technical University
 540 of Denmark, freshmen students received mobile phones in order to study their networks and

541 social behavior in the important change moment of their lives, when joining the university.
542 The study, called SensibleDTU (<https://www.sensible.dtu.dk/?lang=en>), uses not only data
543 collected from the mobile phones (location, Bluetooth-based proximity, call and sms logs etc.)
544 but also data collected from social networks, questionnaires filled out by participants, behavior
545 in economic games and so on. As the data is collected in the context of the university, there is
546 potentially a big issue of students feeling obliged to participate in the study, feeling that their
547 grades may depend on it, or that the data may influence their grades. In this context, we see
548 the implementation of Living Informed Consent not only as a technical mean to put participants
549 in control of the data we collect, but also to convey the message about the opt-in nature of the
550 study, the boundaries of the data usage, and parties accessing the data.

551 It is not feasible to explain the terms and answer all the questions to all 1,000 students
552 personally. The controls must be self-explanatory as much as possible, and guide the user from
553 the first opening of the link to the study to the grant of the authorizations. At the same time,
554 every click made by the user, should be an expression of an informed decision, so the user journey
555 must be a balance of guidance and understanding. For this reason we have created a set of web
556 applications, allowing the users to enroll into the study, express informed consent, and interact
557 with their data.

558 As the study will last for several years, hopefully allowing us to see the life of a student from
559 the very first friendships made until the graduation party, the consent must remain alive. It is
560 again a matter of balance: we do not want the participants to feel under constant surveillance
561 (as they are not, the data is used mostly in aggregated form), at the same time to remember that
562 in fact, the data is being collected and used. We are still trying to understand how to achieve
563 this equilibrium: how often should we remind the users about the collection effort? should they
564 re-authorize applications from time to time? We see a great hope in the applications we create
565 for the users to provide certain services, simple such as life-logging where they can see how
566 active they are, what are their top places etc. and more advanced, such as artistic visualizations
567 of their social networks. Making the user aware of the data by transforming them into value,

568 can greatly benefit the privacy, making users constantly aware what is being collected, but also
569 what kind of value they can get out of it.

570 When a study of such scale is deployed, the particular experiments and sub-studies may
571 not be exactly defined from the very beginning. The initial deployment is a creation of a
572 testbed, where shorter or longer experiments can take place; for example part of the population
573 may participate in the experiment of quantifying the impact of feedback application on their
574 activity levels. Being able to create such experiments in an efficient way is a huge value for the
575 researchers. To do that in the most frictionless way, we give the users the choice to opt-in to
576 those additional experiments, providing some financial or other benefits. This is only possible
577 if there is a notion of identity of the participants, stronger and more useful than a piece of
578 paper with a signature. This identity allows us to reach out to people, offer them additional
579 experiments, and let them agree or disagree to them.

580 This touches upon the re-usability of data, as the new experiments may require additional
581 data to be collected, but also have access to all the existing data, based on user authorization.
582 We can imagine going even further, where entirely different studies can re-use participants data
583 from a previous study based on their authorization. When the data are owned by the users,
584 they are free to authorize access to them to any party that requests it. We can see a New Deal on
585 Data pattern here: rather than services (studies) talking to each other about the user data, they
586 talk directly to the users, seeking their authorization. This can address a very important problem
587 in the research context, the data re-use in a privacy-aware manner. Rather than publishing a
588 static dataset, where the users have lost control over their data, live and fresh data can be
589 continuously accessed by any study that the user agrees to be a part of.

590 Many studies will be willing to offer money or other value for the access to the data. Other
591 will provide the user the opportunity to have new data collected. This way, the data collection
592 becomes an opportunity for the user to enrich their personal dataset, and to benefit from it
593 in the future. Join our study and we will provide you with a smartphone and collect your
594 movement patterns for a year; we will do science and you will gain new data that can get you

595 better value or deals in different services. You may now be eligible for a different study. Or your
 596 music recommendation may get better, because your music service can make a use of this extra
 597 data. Your data.

598 8.2 Scenarios of Use Today, Tomorrow and the Day After

599 **AS: This paragraph is impossible to follow for someone without deep background**
 600 **knowledge of what is the message. Too many random made up scenarios, entities,**
 601 **all mashed together.**

602 By inquiring into and noting the four facets of relevant context described above, it is pos-
 603 sible to describe the basic material contours of any scenario within which Big Data exists such
 604 that the operational framework and adequate approaches to access, use, confidentiality, and
 605 other key interests can be sustainably balanced. In a commercial scenario the relevant people
 606 might be a consumer, merchants, banks, products manufacturers, third party app developers,
 607 and individual members of that consumers bowling team. The relevant transactions might be
 608 a purchase of goods by the consumer from the merchant and the corresponding app that was
 609 embedded in the goods and the downstream transaction of involving the consumer now transact-
 610 ing with the merchant bowling alley and interacting with a bowling team, with whom activity
 611 and sports performance data are shared and aggregated and further mashed up. The rest of
 612 the context can be described for any given scenario and this all could be expressed specifically
 613 rather than by role simply by running a report from the system to indicate it was in fact John
 614 Doe, of openpds.org/owner/571 purchasing a smart bowling ball from Bowl-a-Tronic of [bowlapp-](http://bowlapp-good.com/store/221)
 615 good.com/store/221 and so on for each party that played a role in the relevant scenario. The
 616 same techniques, used for scenarios in other economic sectors and social endeavors shed light
 617 on the fundamental nature and implications of Big Data and options for the use of operational
 618 frameworks acting across domains to balance privacy and access, among other interests.

619 **AS: Bold claims here, not sure if we have sufficient support for them in the**
 620 **chapter.**

621 This book represents a high value opportunity to take stock of the current state and dom-
622 inant trends related to Big Data and help to illuminate important choices at a moment of
623 early adoption, dynamic innovation, and wide open possibilities. By contemplating the relevant
624 contexts of todays scenarios of use in, say, the fields of education, entertainment, government,
625 manufacturing, transportation, and many other core anchors of human activity, we have traction
626 to postulate how todays prevailing trends are likely to result and what changes - perhaps quite
627 small but of profound long term impact - could lead to materially different better outcomes.
628 Consider that if the essence of the New Deal on Data was accepted today, or soon, the na-
629 ture, tenor, capabilities, and experience of living by future generations could be unrecognizingly
630 better. Simply extrapolate from the current anomalous practices regarding personal data and
631 individual identity and push forward the timeline by 5, 10, 20 years and beyond. The current
632 trajectory ends up with dystopian scenarios that effectively reverse hard fought, but easily lost
633 constitutional deal of the United States and social compact of common law societies.

634 By contrast, by adopting the New Deal on Data now it is possible to set conditions that
635 promote prosperity and invention even before the New Deal on Data frameworks are formally
636 launched. This is because the uncertainly and confusion about the basic premises and expecta-
637 tions around personal data and identity will be resolved and so investment and risk taking on
638 a firm foundation can be unleashed. The value of Big Data can be accessed at less direct cost
639 and lower risk when uncertainties about privacy liability are addressed and significant the new
640 value is created by enabling wide scale permission based access to personal data and compu-
641 tations about such data. Adopting use of personal data services in phases, such one economic
642 sector, transaction type or data type at a time enables access to the lower costs and new value
643 in a reasonable manner that allows for time to prepare for and stage each phase of adoption.
644 By staging and phasing the New Deal on Data typical objections to change based on grounds
645 of cost, disruption or over regulation can be addressed. Policy incentives can further address
646 these objections, such as allowing safe harbor protections for conduct of organizations operating
647 under the rules of a trust network. Policy makers can resolve other difficulties by combina-

tions of strategic transition management methods like allowing safe harbor compliance delays, or approving alternative adoption paths and granting other non-substantive waivers to ease any burdens of migrating to new business methods. The key point is change management can be designed to achieve enough value at every phase for every key stakeholder group such that self interests and the broader interests are all aligned with the public good.

9 Future Research (Brian)

Our traditional methods of testing and improving government, organizations, and so on are of limited use in building a data-driven society. Even the scientific method that we normally use do not work as well as we might expect, because there are so many potential connections that our standard statistical tools generate less than useful results.

The reason is that with such rich data, you can easily uncover misleading or unactionable correlations. For instance, let us imagine we discover that people who are unusually active are more likely to get the flu. This is a real example: when we examined the minute-by-minute behavior of a small university community - a real-time flow of gigabytes per day for an entire year - we noticed that an unusual level of running around often predicted onset of the flu [26]. But if we can only analyze the data using traditional statistical methods, we have the problem of discerning why this is true. Is it because the flu virus makes us more active in order to spread itself more quickly? While it is more likely that interacting with many more people than usual makes you more likely to catch the flu, you can't be sure that this is the true cause based on the real-time stream of data alone.

Normal analysis methods do not suffice to answer this type questions, because we do not know all the possible alternatives, and so we cannot form a limited, testable number of clear hypotheses. Instead, we need to devise new ways to test the causality of connections in the real world. We can no longer rely on laboratory experiments; we need to do the experiments in the real world, typically on massive, real-time streams of data.

673 9.1 Research on Design and Deployment of Big Data Systems

674 **AS: I do not understand this paragraph? What is top current research? Where is it**
 675 **applied?** In order to achieve low risk, high value outcomes efficiently, design and deployment
 676 of the coming global wave of Big Data systems should apply top current research. To understand
 677 and address the unique problems and prospects associated with big personal data, the relevant
 678 context must be identified and corresponding rules-driven capabilities must be designed into the
 679 underlying systems.

680 People or systems can determine the right rules to apply to data when the right information
 681 is reliably attached to or logically associated with that data in a standard manner **AS: I think I**
 682 **understand this previous sentences but I' m not sure. What is 'a standard manner'**
 683 **here? What is the right information? It seems it is described in the next sentences,**
 684 **maybe remove this one then?** . Any system that can make, use, receive, or share Big Data
 685 must be capable of associating provenance and purpose for all data in a common and actionable
 686 manner. Requiring a lot of narrative documentation and background about the nuances and
 687 circumstances surrounding every data set is both impractical and counterproductive. By con-
 688 trast, a small amount of metadata listing or reliably linking the parties, transactions, systems
 689 and provenance of the data would suffice. This relevant context together with the data forms
 690 the basis for accountable analysis on big personal data.

691 It is important for science and research to develop further solutions and options ensuring
 692 contextually appropriate rules can be applied by big data systems. For rules to be effectively
 693 applied, systems must not only be able to establish which rules apply but also support the right
 694 functional capabilities and have appropriate information structure, format, and meta-data.

695 Some capabilities will likely be essential to all Big Data systems, such as highly scalable
 696 active storage, standard methods for integration with other Big Data systems, and a processing
 697 architecture enabling high speed statistical analytics. But there are and will continue to emerge
 698 multiple types of Big Data systems. Some functions or controls will likely be important —
 699 or even feasible — only for certain types of future systems. For instance, it is reasonable to

700 expect some systems will specialize in enormous volumes of entirely non-personal data from
 701 many real-time sources (e.g. for soil science, materials engineering, astronomy) while other Big
 702 Data systems will hinge upon mass quantities of highly sensitive personal information (e.g. for
 703 clinical medicine, education and life-long learning, social entertainment).

704 **AS: I feel Big Data term is abused in this section...**

705 While some capabilities, such as ingesting and processing astronomical data-sets, will be
 706 unique to only a subset of Big Data systems, it is reasonable to anticipate that data will be
 707 increasingly cross-tabulated, merged, and otherwise shared with other systems and data. It can
 708 be nearly impossible to conclusively predict for the entire life of a system what data will be
 709 received by, created in, or transmitted from that system at the design phase. This prediction is
 710 all the harder to make when the systems are intended for Big Data.

711 The four contextual facets of people, interactions, technology, and data provide a sound
 712 underpinning for the design of new Big Data and Web 2.0 systems. The existing systems design
 713 and development processes of establishing business cases, use cases, agile stories, functional
 714 requirements, etc. do not reliably identify the factors most relevant to use of Big Data, especially
 715 in a Web 2.0 massively distributed environment. The four facets can also be used to analyze
 716 appropriate, required or prohibited uses for existing Big Data systems. However, it can be
 717 difficult to extract the relevant information from or apply any effective control on systems used
 718 for Big Data but designed to achieve limited purposes in hierarchical closed environments.

719 Big Data, by its nature, represents a new set of business, legal, and technical capabilities and
 720 requirements. Most of the worlds systems today are not capable of ingesting, storing, using, or
 721 dynamically flowing big data with other systems. Considering that a) Big Data is of high value
 722 immediately and higher value in the short and long terms, and b) the young but competitive
 723 marketplace of Big Data system components, platforms, applications, and other solutions is a
 724 hotbed of innovation it can be predicted that a transition to Big Data systems will continue.
 725 The key observation is that virtually all Big Data systems have yet to be designed, implemented,
 726 customized, or deployed. Institutions that are the current early adopters of todays Big Data

727 system will soon replace those systems and the rest of the world will adopt big data systems in
 728 phases over time. Based upon this observation, **AS: ??????????????**

729 **9.2 Research on Big Data for Design of Institutions**

730 Using massive, live data to design institutions and policies is outside of our normal way of
 731 managing things. We live in an era that builds on centuries of science and engineering, and
 732 the standard choices for improving systems, governments, organizations, and so on are fairly
 733 well understood. Therefore our scientific experiments normally need only consider a few clear
 734 alternatives, ‘plausible hypotheses’.

735 With the coming of Big Data, we are going to be operating very much out of our old,
 736 familiar ballpark. These data are often indirect and noisy, and so interpretation of the data
 737 requires greater care than usual. Even more importantly, a great deal of the data is about
 738 human behavior, and the questions are ones that seek to connect physical conditions to social
 739 outcomes. Until we have a solid, well-proven, and quantitative theory of social physics, we will
 740 not be able to formulate and test hypotheses in the way we can when we design bridges or
 741 develop new drugs.

742 Therefore, we must move beyond the closed, laboratory-based question-and-answering pro-
 743 cess that we currently use, and begin to manage our society in a new way. We must begin to test
 744 connections in the real world far earlier and more frequently than we have ever had to do before,
 745 using the methods the Human Dynamics research group have developed with our collaborators
 746 for the Friends and Family [3] or the SensibleDTU (<https://www.sensible.dtu.dk>) study. We
 747 need to construct Living Laboratories — communities willing to try a new way of doing things
 748 or, to put it bluntly, to be guinea pigs — in order to test and prove our ideas. This is new
 749 territory and so it is important for us to constantly try out new ideas in the real world in order
 750 to see what works and what does not.

751 An example of such a Living Lab is the ‘open data city just launched by one author (Pentland)
 752 with the city of Trento in Italy, along with Telecom Italia, Telefonica, the research university

Fondazione Bruno Kessler, the Institute for Data Driven Design, and local companies. Importantly, this Living Lab has the approval and informed consent of all its participants they know that they are part of a gigantic experiment whose goal is to invent a better way of living. More detail on this Living Lab can be found at <http://www.mobileterritoriallab.eu/>.

The goal of this Living Lab is to develop new ways of sharing data to promote greater civic engagement and exploration. One specific goal is to build upon and test trust-network software such as our openPDS system. Tools such as openPDS make it safe for individuals to share personal data (e.g., health data, facts about your children) by controlling where your data go and what is done with them.

The specific research questions we are exploring depend upon a set of “personal data services” designed to enable users to collect, store, manage, disclose, share, and use data about themselves. These data can be used for the personal self-empowerment of each member, or (when aggregated) for the improvement of the community through data commons that enable social network incentives. The ability to share data safely should enable better idea flow among individuals, companies, and government, and we want to see if these tools can in fact increase productivity and creative output at the scale of an entire city.

An example of an application enabled by the openPDS trust frame work is sharing of best practices among families with young children. How do other families spend their money? How much do they get out and socialize? Which preschools or doctors do people stay with for the longest time? Once the individual gives permission, our openPDS system allows such personal data to be collected, anonymized, and shared with other young families safely and automatically.

The openPDS system lets the community of young families learn from each other without the work of entering data by hand or the risk of sharing through current social media. While the Trento experiment is still in its early days, the initial reaction from participating families is that these sorts of data sharing capabilities are valuable, and they feel safe sharing their data using the openPDS system.

The Trento Living Lab will let us investigate how to deal with the sensitivities of collecting

and using deeply personal data in real-world situations. In particular, the Lab will be used as a pilot for the New Deal on Data and for new ways to give users control of the use of their personal data. For example, we will explore different techniques and methodologies to protect the users privacy while at the same time being able to use these personal data to generate a useful data commons. We will also explore different user interfaces for privacy settings, for configuring the data collected, for the data disclosed to applications and for those shared with other users, all in the context of a trust framework.

10 Conclusions

Our societies today face unprecedented challenges. Solving those problems will require access to the personal data, so we can understand how the society works, how we move around, what makes us productive, how the ideas and diseases spread. The insights must be actionable, available in real-time, and engaging the population, creating the nervous system of the society. In this chapter we have reviewed how Big Data collected in institutional context can be used for the public good. In many cases, the data needed for creating better society is already collected and exists closed in silos of companies and governments. Using well designed and implemented set of institutional controls, covering business, legal, and technical dimensions, we described how the silos can be opened. The framework for doing this — the New Deal on Data — postulates that the primary driver of the change must be the ownership of the personal data, given to people about whom the data is. This ownership, the right to use, transfer, and remove the data ensures that the data is available for public good, while at the same time protecting the privacy of the citizens.

The New Deal on Data is still new. Here we described our efforts in understanding the technical means of how it can be implemented, the legal framework around it, business ramifications, and the direct value that can be derived from researchers, companies, governments, and users having more access to the data. It is clear that companies must play the major role in the implementation of the New Deal, incentivized by business opportunities and pressured by the

806 legislation and demand of the users. Only with such orchestration it will be possible to change
 807 the current feudal system of the data ownership and finally put the immense quantities of the
 808 collected personal data to good use.

809 References

- 810 1. Binding obligations on User-Managed Access (UMA) participants. Technical Specifica-
 811 tions draft-maler-oauth-umatrust-01, Kantara Initiative, July 2013.
- 812 2. User-Managed Access (UMA) profile of OAuth2.0. Technical Specifications draft-
 813 hardjono-oauth-umacore-08, Kantara Initiative, December 2013.
- 814 3. Nadav Aharony, Wei Pan, Cory Ip, Inas Khayal, and Alex Pentland. Social fmri: Investi-
 815 gating and shaping social mechanisms in the real world. *Pervasive and Mobile Computing*,
 816 7(6):643–659, 2011.
- 817 4. Sinan Aral and Dylan Walker. Identifying influential and susceptible members of social
 818 networks. *Science*, 337(6092):337–341, 2012.
- 819 5. Simon Barber, Xavier Boyen, Elaine Shi, and Ersin Uzun. Bitter to Better – how to
 820 make Bitcoin a better currency. In *Proceedings Financial Cryptography and Data Security*
 821 *Conference (Lecture Notes in Computer Science Volume 7397)*, pages 399–414, April 2012.
- 822 6. Ellen Barry. Protests in moldova explode, with help of twitter. *New York Times*, 8, 2009.
- 823 7. Nick Bilton. Girls around me: An app takes creepy to a new level. *The New York Times*.
- 824 8. Center for Environmental & Human Toxicology University of Florida. Development of
 825 Cleanup Target Levels (CTLs) For Chapter 62-777, F.A.C. Technical report, Division of
 826 Waste Management Florida Department of Environmental Protection, February 2005.

- 827 9. Paul Lukowicz Bert Arnrich Cornelia Setz Gerhard Troster David Tacconi, Oscar Mayora
828 and Christian Haring. Activity and emotion recognition to support early diagnosis of
829 psychiatric diseases. pages 100–102. IEEE, 2008.
- 830 10. Yves-Alexandre de Montjoye, César A Hidalgo, Michel Verleysen, and Vincent D Blondel.
831 Unique in the crowd: The privacy bounds of human mobility. *Scientific reports*, 3, 2013.
- 832 11. Yves-Alexandre de Montjoye, Samuel S Wang, Alex Pentland, Dinh Tien Tuan Anh, An-
833 witaman Datta, Kevin W Hamlen, Lalana Kagal, Murat Kantarcioglu, Vaibhav Khadilkar,
834 Kerim Yasin Oktay, et al. On the trusted use of large-scale personal data. *IEEE Data*
835 *Eng. Bull.*, 35(4):5–8, 2012.
- 836 12. Ralph A. DeMeo and Sarah Meyer Doar. Restrictive covenants as institutional controls
837 for remediated sites: Worth the effort? *The Florida Bar Journal*, 85(2), 2011.
- 838 13. Nathan Eagle and Alex Pentland. Reality mining: sensing complex social systems. *Per-*
839 *sonal and ubiquitous computing*, 10(4):255–268, 2006.
- 840 14. Jonathan Woetzel et al. Preparing for china’s urban billion. 2009.
- 841 15. Florida Department of Environmental Protection - Division of Waste Management. Insti-
842 tutional Controls Procedures Guidance. [http://www.dep.state.fl.us/waste/quick\](http://www.dep.state.fl.us/waste/quick_topics/publications/wc/csf/icpg.pdf)
843 [_topics/publications/wc/csf/icpg.pdf](http://www.dep.state.fl.us/waste/quick_topics/publications/wc/csf/icpg.pdf), June 2012.
- 844 16. Kim Gittleson. How big data is changing the cost of insurance. *BBC News*, 2013.
- 845 17. Kate Greene. Reality mining. *Technology Review*, 2008.
- 846 18. Lev Grossman. Iran protests: Twitter, the medium of the movement. *Time Magazine*,
847 17, 2009.
- 848 19. Aniko Hannak, Piotr Sapiezynski, Arash Molavi Kakhki, Balachander Krishnamurthy,
849 David Lazer, Alan Mislove, and Christo Wilson. Measuring personalization of web search.

- 850 In *Proceedings of the 22nd international conference on World Wide Web*, pages 527–538.
 851 International World Wide Web Conferences Steering Committee, 2013.
- 852 20. Thomas Hardjono, Patrick Deegan, and John Clippinger. On the Design of Trustworthy
 853 Compute Frameworks for Self-Organizing Digital Institutions. In *Proceedings of the 16th*
 854 *International Conference on Human-Computer Interaction*, 2014.
- 855 21. Thomas Hardjono, Daniel Greenwood, and Alex Pentland. Towards a trustworthy digital
 856 infrastructure for core identities and personal data stores. In *Proceedings of the ID360*
 857 *Conference on Identity*. University of Texas, April 2013.
- 858 22. Juniper Networks. Secure Data Access Anywhere and Anytime: Current Landscape and
 859 Future Outlook of Enterprise Mobile Security. A forrester consulting thought leadership
 860 paper commissioned by att and juniper networks, Forrester Research, October 2012.
- 861 23. Meglena Kuneva. Roundtable on Online Data Collection, Targeting and Profiling . [http:](http://europa.eu/rapid/press-release_SPEECH-09-156_en.htm)
 862 [//europa.eu/rapid/press-release_SPEECH-09-156_en.htm](http://europa.eu/rapid/press-release_SPEECH-09-156_en.htm), 2009.
- 863 24. David Lazer, Alex Sandy Pentland, Lada Adamic, Sinan Aral, Albert Laszlo Barabasi,
 864 Devon Brewer, Nicholas Christakis, Noshir Contractor, James Fowler, Myron Gutmann,
 865 et al. Life in the network: the coming age of computational social science. *Science (New*
 866 *York, NY)*, 323(5915):721, 2009.
- 867 25. Antonio Lima, Manlio De Domenico, Veljko Pejovic, and Mirco Musolesi. Exploiting
 868 cellular data for disease containment and information campaigns strategies in country-
 869 wide epidemics. School of computer science university of birmingham technical report
 870 csr-13-01, University of Birmingham, May 2013.
- 871 26. Anmol Madan, Manuel Cebrian, David Lazer, and Alex Pentland. Social sensing for
 872 epidemiological behavior change. In *Proceedings of the 12th ACM international conference*
 873 *on Ubiquitous computing*, pages 291–300. ACM, 2010.

- 874 27. AC Madrigal. Dark social: We have the whole history of the web wrong. *The Atlantic*,
875 2013.
- 876 28. Alan Mislove, Sune Lehmann, Yong-Yeol Ahn, Jukka-Pekka Onnela, and J Niels Rosen-
877 quist. Pulse of the nation: Us mood throughout the day inferred from twitter. *Accessed*
878 *November*, 22(2011):2011, 2010.
- 879 29. Arvind Narayanan and Vitaly Shmatikov. Robust de-anonymization of large sparse
880 datasets. In *Security and Privacy, 2008. SP 2008. IEEE Symposium on*, pages 111–125.
881 IEEE, 2008.
- 882 30. Wei Pan, Yaniv Altshuler, and Alex Sandy Pentland. Decoding social influence and
883 the wisdom of the crowd in financial trading network. In *Privacy, Security, Risk and*
884 *Trust (PASSAT), 2012 International Conference on and 2012 International Confernece*
885 *on Social Computing (SocialCom)*, pages 203–209. IEEE, 2012.
- 886 31. Wei Pan, Gourab Ghoshal, Coco Krumme, Manuel Cebrian, and Alex Pentland. Urban
887 characteristics attributable to density-driven tie formation. *Nature communications*, 4,
888 2013.
- 889 32. ALEX PENTLAND. Reality mining of mobile communications: Toward a new deal on
890 data. *The Global Information Technology Report 2008–2009*, page 1981, 2009.
- 891 33. Alex Pentland, David Lazer, Devon Brewer, and Tracy Heibeck. Using reality mining to
892 improve public health and medicine. *Stud Health Technol Inform*, 149:93–102, 2009.
- 893 34. Vivek K Singh, Laura Freeman, Bruno Lepri, and Alex Sandy Pentland. Classifying
894 spending behavior using socio-mobile data. *HUMAN*, 2(2):pp–99, 2013.
- 895 35. Chaoming Song, Zehui Qu, Nicholas Blumm, and Albert-László Barabási. Limits of
896 predictability in human mobility. *Science*, 327(5968):1018–1021, 2010.
- 897 36. Stan Stalnaker. The Ven currency, 2013. <http://www.ven.vc>.

- 898 37. Latanya Sweeney. Simple demographics often identify people uniquely. *Health (San Fran-*
899 *cisco)*, pages 1–34, 2000.
- 900 38. The White House. National Strategy for Trusted Identities in Cyberspace: Enhancing On-
901 line Choice, Efficiency, Security, and Privacy. The White House, April 2011. Available on
902 http://www.whitehouse.gov/sites/default/files/rss_viewer/NSTICstrategy_041511.pdf.
- 903 39. United States Environmental Protection Agency. Institutional Controls Bibliography.
904 <http://www.epa.gov/superfund/policy/ic/guide/biblio.pdf>, December 2005.
- 905 40. United States Environmental Protection Agency. RCRA Corrective Action Institu-
906 tional Controls - glossary. [http://www.epa.gov/epawaste/hazard/correctiveaction/](http://www.epa.gov/epawaste/hazard/correctiveaction/resources/guidance/ics/glossary1.pdf)
907 [resources/guidance/ics/glossary1.pdf](http://www.epa.gov/epawaste/hazard/correctiveaction/resources/guidance/ics/glossary1.pdf), 2007.
- 908 41. United States Environmental Protection Agency. Institutional Controls: A Guide to Plan-
909 ning, Implementing, Maintaining, and Enforcing Institutional Controls at Contaminated
910 Sites. Technical Report OSWER 9355.0-89 EPA-540-R-09-001, EPA, December 2012.
- 911 42. Jessica Vitak, Paul Zube, Andrew Smock, Caleb T Carr, Nicole Ellison, and Cliff Lampe.
912 It’s complicated: Facebook users’ political participation in the 2008 election. *CyberPsy-*
913 *chology, behavior, and social networking*, 14(3):107–114, 2011.
- 914 43. World Economic Forum. Personal Data: The Emergence of a New
915 Asset Class, 2011. Available on [http://www.weforum.org/reports/](http://www.weforum.org/reports/personal-data-emergence-new-asset-class)
916 [personal-data-emergence-new-asset-class](http://www.weforum.org/reports/personal-data-emergence-new-asset-class).