

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	2
10	2 The New Deal on Data	4
11	3 Personal Data: Emergence of a New Asset Class	6
12	4 Enforcing the New Deal on Data	10
13	5 Transitioning End-User Assent Practices	13
14	6 Business, Legal, and Technical Dimensions of Big Data Systems	15
15	7 Big Data and Personal Data Institutional Controls	16
16	8 Scenarios of Use in Context	21
17	8.1 Example Scenario: Research Systems	21
18	8.2 Scenarios of Use Today, Tomorrow and the Day After	24
19	9 Future Research	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

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 the New Deal on Data, and describe
29 future directions for 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 controls must support and reflect
34 greater user control over personal data, as well as large scale interoperability for data sharing be-
35 tween and among institutions. Core capabilities of these controls include responsive rule-based
36 systems governance and fine-grained authorizations for distributed rights management.

37 Sustaining a healthy, safe, and efficient society is a scientific and engineering challenge going
38 back to the 1800s when the Industrial Revolution spurred rapid urban growth, thereby 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 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 energy efficient, have secure food
49 and water supplies, are protected from pandemics and enjoy much better governance. To reach
50 these 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. Finally, we need to create channels for society to agree upon and communicate those
55 needs.

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

63 The engine driving this nervous system is Big Data: the newly ubiquitous digital data, now
64 available about all aspects of human life. We can analyze patterns of human experience and
65 ideas exchange within the *digital breadcrumbs* that we all leave behind as we move through the
66 world: call records, credit card transactions, GPS location fixes, among others. By recording
67 our choices, these data tell the story of our lives. And this may be very different from what
68 we decide to put on Facebook or Twitter; our postings there are what we choose to tell people,
69 edited according to the standards of the day and filtered to match the persona we are building.
70 Mining social networks can give some great insights about human nature [4, 27, 41]; who we
71 really are, however, is even more accurately determined by where we spend our time and which
72 things we buy, rather than just what we say we do [26].

73 The process of analyzing the patterns within these digital breadcrumbs is called reality
74 mining [13, 31], and through it we can learn an enormous amount about who we are. The

75 Human Dynamics research group at MIT found that we can use them to tell if we are likely
 76 to get diabetes [32], or whether we are the sort of person who will pay back loans [33]. By
 77 analyzing these patterns across many people, we are discovering that we can begin to explain
 78 many things — crashes, revolutions, bubbles — that previously appeared to be random acts of
 79 God [29]. For this reason, the magazine Technology Review named our development of reality
 80 mining as one of the ten technologies that will change the world [16].

81 2 The New Deal on Data

82 The digital breadcrumbs we leave behind provide clues about who we are, what we do and what
 83 we want. This makes personal data — data about individuals — immensely valuable, both for
 84 public good and for private companies. As European Consumer Commissioner, Meglena Kuneva
 85 said recently, “Personal data is the new oil of the Internet and the new currency of the digital
 86 world” [22]. This new ability to see the details of every interaction can be used for good or for
 87 ill. Therefore, maintaining protection of personal privacy and freedom is critical to our future
 88 success as a society. We need to enable even more data sharing for the public good; at the same
 89 time, we need to do a much better job in protecting the privacy of the individuals.

90 A successful data-driven society must be able to guarantee that our data will not be abused;
 91 perhaps especially that government will not abuse the power conferred by access to such fine-
 92 grain data. The abuses may be directly targeted at users, for example by offering them higher
 93 insurance rates based on their shopping history, or create problems for the entire society in the
 94 long run, for example by limiting user choices and closing them into information bubbles [18]. To
 95 achieve the positive possibilities of the new society, we require the *New Deal on Data*, workable
 96 guarantees that the data needed for public good are readily available while at the same time
 97 protecting the citizenry [31].

98 The key insight that motivates the idea of the New Deal on Data is that our data are worth
 99 more when shared, because these aggregated data — averaged, combined across population, and
 100 often distilled to high-level features — inform improvements in systems such as public health,

101 transportation, and government. For instance, we have demonstrated that data about the way
102 we behave and where we go can be used to minimize the spread of infectious disease [25,32]. Our
103 research has reported how we were able to use these digital breadcrumbs to track the spread of
104 influenza from person to person on an individual level. And if we can see it, we can stop it.

105 Similarly, if we are worried about global warming, these shared, aggregated data can show us
106 how patterns of mobility relate to productivity [30]. In turn, this provides us with the ability to
107 design cities that are more productive and, at the same time, more energy efficient. But in order
108 to obtain these results and make a greener world, we need to be able to see the people moving
109 around; this depends on many people willing to contribute their data, even if only anonymously
110 and in aggregate.

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

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

128 The ultimate goal is to provide the society with tools to analyze and understand what needs
 129 to be done, and to reach the consensus on how to do it. This goes beyond just creating more
 130 communication platforms; the assumption that more interactions between users will result in
 131 better decisions being made may be very misleading. Although in the recent years we have seen
 132 some great examples of using social networks for better organization in society, for example
 133 during political protests [6, 17], we are not even close to the point where we can start reaching
 134 consensus about the big problems: epidemics, climate change, pollution. The discussions must
 135 be data driven, involving both experts and wisdom of the crowds – users themselves interested
 136 in improving the society. The problems we are dealing with now as a global society are not
 137 easy. We are responsible for many of them, and being able to tackle them on a global scale is
 138 necessary for our survival as a people.

139 **3 Personal Data: Emergence of a New Asset Class**

140 It has long been recognized that the first step to promoting liquidity in land and commodity mar-
 141 kets is to guarantee ownership rights so that people can safely buy and sell. Similarly, the first
 142 step toward creating more new ideas and greater flow ideas (aka idea liquidity) is to define own-
 143 ership rights. The only politically viable course is to give individual citizens key rights over data
 144 that are about them and in fact, these types of rights have undergirded the European Union’s
 145 Privacy Directive since 1995 (See: [http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?](http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=CELEX:31995L0046:EN:HTML)
 146 [uri=CELEX:31995L0046:EN:HTML](http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=CELEX:31995L0046:EN:HTML))

147 We need to recognize personal data as a valuable asset of the individual that is given to
 148 companies and government in return for services.

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

- 151 • You have the right to possess data about you. Regardless of what entity collects the data,
 152 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 land ownership disputes, for example.

168 In 2007, one author (Pentland) first proposed the New Deal on Data to the World Economic
 169 Forum [42]. 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 the public good, while
 173 at 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 [42]. 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)
183 and the citizen (the user-base of these services) allowed the developed nations to expand their
184 economies 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 [42], 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 enroll and register end-users. These siloed
194 repositories of personal data exemplify 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
202 into the advertising industry. Personal privacy concerns are thus inadequately addressed at
203 best, or simply overlooked in the majority of 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 [24].
 213 In releasing and analyzing this data, the privacy of the people who generated the data was
 214 protected not only by technical means, such as removal of 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 [28] and other similar privacy breaches [36], 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 [34], we are also very unique. Having access to one dataset
 220 may be enough to uniquely fingerprint someone based on just a 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 [42] 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 [37] — albeit still in its infancy —
 227 represent 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 [32], 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

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 upon 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 for a given trust network arise from that systems underlying trust models and the contextual scenarios within which the networked data and the relationships of parties occur.

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 multi-party 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

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

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

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

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

308 The Human Dynamics Lab has applied the system rules approach to development of in-
309 tegrated business, technical architecture, and rules large scale institutional use of personal
310 data stores, available as an example under MIT’s creative commons license by MIT, at [https:](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. Therefore,
 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**

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

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

337 promote the privacy of the users, effectively pushing them to press *I Agree* in every presented
338 window.

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

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

360 6 Business, Legal, and Technical Dimensions of Big Data Sys- 361 tems

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

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

382 Defining as a ‘system’ the thing to which institutional controls apply provides an achievable
383 and measurable basis for balancing privacy, access and other interests in Big Data. Within a
384 given institution, there may in fact be many different discernable organizations and correspond-
385 ing systems. Meanwhile the system of one institution frequently exists across many different

external institutions. The application of Big Data institutional controls can be applied across the board to a unit of a given institution or targeted by agreement to certain types of data or particular transactions spanning many institutions. Once a systems view is adopted, there is a tractable starting point to narrow or broaden the scope of view, to focus on material dimensions of a system and therefore enable more effective use and control of Big Data.

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

7 Big Data and Personal Data Institutional Controls

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

A good example of how this concept supports and reflects the goals and objectives of environmental regulation can be found in the policy documents of the Environmental Protection Agency (EPA). This following definition is instructive, and is part of the Institutional Control Glossary of Terms [39]:

“Institutional Controls - Non-engineering measures intended to affect human activities in such a way as to prevent or reduce exposure to hazardous substances. They are almost always used in conjunction with, or as a supplement to, other measures such as waste treatment or containment. There are four categories of institutional

411 controls: governmental controls; proprietary controls; enforcement tools; and infor-
412 mational devices.”

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

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

419 In domains of information technology, this approach is most commonly reflected as “enter-
420 prise controls” related to security. See, for example, the Juniper Networks enterprise security
421 report [21] stating: “Enterprise mobility technologies, especially those designed to retrofit en-
422 terprise controls on top of consumer mobile devices, are rapidly evolving. This was a message
423 we heard loud and clear in the study.” This study and analysis also reveals much about the
424 internal controls needed to accommodate mobile device use by employees. In both capacities as
425 employee, consumer, and other roles, the use of mobile devices triggers myriad legal, policy, and
426 other implications for institutional controls.

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

437 Consider again the analogy of institutional controls in the context of environmental law, and
438 how these types of measures can be applied in the Big Data, privacy, and access context to
439 digital environments. Given the relatively mature and stable state of environmental regulation,
440 there is much to be learned by examining this context of institutional controls. Environmental
441 regulatory compliance with waste management cleanup requirements could include institutional
442 controls restricting land use on adjacent property. In these situations, it is possible that the
443 remediation strategy requires significant use of land outside the property boundaries of the
444 cleanup site. In these cases, the regulators and the land owner responsible for the regulated
445 property must find ways to ensure a common approach among multiple owners and across
446 multiple property environments. Clauses on the relevant deeds, an enforceable consent order,
447 or targeted regulations and zoning rules are examples of more severe institutional controls that
448 can be employed to ensure consistent and effective actions are taken across ownership and real
449 property boundaries.

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

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

463 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 Controls* [39]. Also see: *Institutional Controls Bibliography: Institutional Control, Remedy Selection, and Post-Construction Completion Guidance and Policy, December 2005* [38].

When institutional controls would apply to “separately owned neighboring properties” a number of issues arise that are very relevant to the problems associated with managing personal and big data across legal, business and other systemic boundaries. Requiring the party responsible for site cleanup to use “best efforts” to attain agreement by third parties to institute the relevant institutional controls is perhaps the most direct and least prescriptive approach. When direct negotiated agreement is not successful, then use of third party neutrals to resolve disagreements regarding institutional controls can be required. If necessary, environmental regulation can force an acquisition of neighboring land by compelling the party responsible to purchase the other property or by purchase of the property directly by the EPA [40].

In the context of Big Data, privacy, and access, institutional controls are seldom, if ever, the result of government regulatory frameworks such as are seen in the environmental waste management oversight by the EPA. Rather, institutions applying measures constituting institutional controls in the Big Data and related information technology and enterprise architecture contexts will typically employ governance safeguards, business practices, legal contracts, technical security, reporting, and audit programs and various risk management measures.

Inevitably, institutional controls for Big Data will have to operate effectively across institutional boundaries, just as environmental waste management internal controls must sometimes be applied across real property boundaries and may subject multiple different owners to enforcement actions corresponding to the applicable controls. Short of government regulation, the use of system rules as a general model are one widely understood, accepted, and efficient method for defining, agreeing, and enforcing institutional and other controls across business, legal, and technical domains of ownership, governance, and operation.

The use of system rules and integrated participation agreements by developers and end-

491 users is a way to ensure intended operational frameworks conform to applicable institutional
 492 controls. The example of Living Informed Consent described in this chapter, demonstrates how
 493 institutional controls comprised of legal and definite workflow measures, in concert with technical
 494 methods, can result in a higher level of performance, while appropriately balancing legitimate
 495 interests of various parties regarding use and access to personal data.

496 Following the World Economic Forum recommendations of treating personal data stores in
 497 the manner of bank accounts [42], there are a number of infrastructure improvements that need to
 498 be realized, if the personal data ecosystem is to flourish and deliver new economic opportunities.
 499 We believe the following infrastructure improvements are necessary for the coming personal data
 500 ecosystem:

- 501 • *New global data provenance network:* In order for personal data to be treated like bank
 502 accounts, the origin information regarding data items coming into the data store must be
 503 maintained [20]. In other words, the provenance of all data items must be accounted for
 504 by the IT infrastructure upon which the personal data store operates. The heterogeneous
 505 provenance databases must then be interconnected in order to provide a resilient and
 506 scalable platform for audit and accounting systems to track and reconcile the movement
 507 of personal data from the respective data stores.
- 508 • *Trust network for computational law:* In order for trust to be established between parties
 509 who wish to exchange personal data, we foresee that some degree of “computational law”
 510 technologies may have to be integrated into the design of personal data systems. Such
 511 technologies should not only verify terms of contracts (e.g. terms of data use) against user-
 512 defined policies but also have mechanisms built-in to ensure non-repudiation of entities who
 513 have accepted these digital contracts. Efforts such as [1,2] are beginning to bring better
 514 evidentiary proof and enforceability of contracts into the technical protocol flows.
- 515 • *Development of institutional controls for digital institutions:* Currently there are a number
 516 of proposals for the creation of virtual currencies (e.g. BitCoin [5], Ven [35]) in which the

517 systems have the potential to evolve into self-governing “digital institutions” [19]. Such
 518 systems and institutions that operate on them will necessitate the development of a new
 519 paradigm to understand the aspects of institutional control within their context.

520 8 Scenarios of Use in Context

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

528 8.1 Example Scenario: Research Systems

529 Computational Social Science (CSS) studies are based on data collected often with an extremely
 530 high resolution and scale [23]. Using computational power combined with mathematical models,
 531 such data can be used to provide insights into human nature. Much of the data collected, for
 532 example mobility traces are sensitive and private; most individuals would feel uncomfortable
 533 sharing them publicly. The need for solutions to ensure the privacy of the individuals has grown
 534 alongside the data collection efforts.

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

542 that participants in fact expressed informed consent.

543 One author (Stopczynski) deployed this year a 1,000 phones study at Technical University
544 of Denmark, freshmen students received mobile phones in order to study their networks and
545 social behavior in the important change moment of their lives, when joining the university.
546 The study, called SensibleDTU (<https://www.sensible.dtu.dk/?lang=en>), uses not only data
547 collected from the mobile phones (location, Bluetooth-based proximity, call and sms logs etc.)
548 but also data collected from social networks, questionnaires filled out by participants, behavior
549 in economic games and so on. As the data is collected in the context of the university, there is
550 potentially a big issue of students feeling obliged to participate in the study, feeling that their
551 grades may depend on it, or that the data may influence their grades. In this context, we see
552 the implementation of Living Informed Consent not only as a technical mean to put participants
553 in control of the data we collect, but also to convey the message about the opt-in nature of the
554 study, the boundaries of the data usage, and parties accessing the data.

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

562 As the study will last for several years, hopefully allowing us to see the life of a student from
563 the very first friendships made until the graduation party, the consent must remain alive. It is
564 again a matter of balance: we do not want the participants to feel under constant surveillance
565 (as they are not, the data is used mostly in aggregated form), at the same time to remember that
566 in fact, the data is being collected and used. We are still trying to understand how to achieve
567 this equilibrium: how often should we remind the users about the collection effort? should they
568 re-authorize applications from time to time? We see a great hope in the applications we create

569 for the users to provide certain services, simple such as life-logging where they can see how
570 active they are, what are their top places etc. and more advanced, such as artistic visualizations
571 of their social networks. Making the user aware of the data by transforming them into value,
572 can greatly benefit the privacy, making users constantly aware what is being collected, but also
573 what kind of value they can get out of it.

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

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

594 Many studies will be willing to offer money or other value for the access to the data. Other
595 will provide the user the opportunity to have new data collected. This way, the data collection

becomes an opportunity for the user to enrich their personal dataset, and to benefit from it in the future. Join our study and we will provide you with a smartphone and collect your movement patterns for a year; we will do science and you will gain new data that can get you better value or deals in different services. You may now be eligible for a different study. Or your music recommendation may get better, because your music service can make a use of this extra data. Your data.

8.2 Scenarios of Use Today, Tomorrow and the Day After

By inquiring into and noting the four facets of relevant context described above, it is possible to describe the basic material contours of any scenario within which Big Data exists such that the operational framework and adequate approaches to access, use, confidentiality, and other key interests can be sustainably balanced. In a commercial scenario the relevant people might be a consumer, merchants, banks, products manufacturers, third party app developers, and individual members of that consumers bowling team.

The relevant transactions might be a purchase of goods by the consumer from the merchant and the corresponding app that was embedded in the goods and the downstream transaction of involving the consumer now transacting with the merchant bowling alley and interacting with a bowling team, with whom activity and sports performance data are shared and aggregated and further mashed up. The rest of the context can be described for any given scenario and this all could be expressed specifically rather than by role simply by running a report from the system to indicate it was in fact John Doe, of openpds.org/owner/571 purchasing a smart bowling ball from Bowl-a-Tronic of bowlappgood.com/store/221 and so on for each party that played a role in the relevant scenario. The same techniques, used for scenarios in other economic sectors and social endeavors shed light on the fundamental nature and implications of Big Data and options for the use of operational frameworks acting across domains to balance privacy and access, among other interests.

The New Deal on Data is designed to provide good value to all stakeholders creating, using

622 or benefiting from personal data, but the entire vision need not be adopted before value starts to
623 flow. The social science research study scenario (below) demonstrates how researchers and study
624 participants alike derive value from New Deal on Data principles today. As more researchers
625 and students use these types of systems the value is predicted to increase based upon a network
626 effect.

627 Adopting New Deal on Data principles on a large scale can be accomplished iteratively, such
628 one economic sector, transaction type or data type at a time. A reasonable success metric for
629 adoption of large scale visions such as the New Deal on Data is whether change management
630 has been designed to achieve enough value at every phase for every key stakeholder group to
631 make the change worth the effort. Value to all parties participating in the New Deal on Data
632 increases as direct or indirect use and re-use of personal data is available in greater volumes and
633 varieties. Such volume and variety of personal data increases as more parties and transaction
634 types and data sets and systems adopt and interoperate within the New Deal on Data.

635 By staging and phasing adoption of the New Deal on Data typical objections to change based
636 on grounds of cost, disruption or over regulation can be addressed. Policy incentives can further
637 address these objections, such as allowing safe harbor protections for conduct of organizations
638 operating under the rules of a trust network. Policy makers can resolve other difficulties by
639 combinations of strategic transition management methods like allowing safe harbor compliance
640 delays, or approving alternative adoption paths and granting other non-substantive waivers to
641 ease any burdens of migrating to new business methods.

642 9 Future Research

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

647 The reason is that with such rich data, you can easily uncover misleading or unactionable

648 correlations. For instance, let us imagine we discover that people who are unusually active are
 649 more likely to get the flu. This is a real example: when we examined the minute-by-minute
 650 behavior of a small university community - a real-time flow of gigabytes per day for an entire
 651 year - we noticed that an unusual level of running around often predicted onset of the flu [25].
 652 But if we can only analyze the data using traditional statistical methods, we have the problem
 653 of discerning why this is true. Is it because the flu virus makes us more active in order to spread
 654 itself more quickly? While it is more likely that interacting with many more people than usual
 655 makes you more likely to catch the flu, you can't be sure that this is the true cause based on
 656 the real-time stream of data alone.

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

662 9.1 Research on Design and Deployment of Big Data Systems

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

669 People or systems can determine the right rules to apply to data when the right information
 670 is reliably attached to or logically associated with that data in a standard manner **AS: I think**
 671 **I understand this previous sentences but I'm not sure. What is 'a standard manner'**
 672 **here? What is the right information? It seems it is described in the next sentences,**
 673 **maybe remove this one then?** . Any system that can make, use, receive, or share Big Data

674 must be capable of associating provenance and purpose for all data in a common and actionable
 675 manner. Requiring a lot of narrative documentation and background about the nuances and
 676 circumstances surrounding every data set is both impractical and counterproductive. By con-
 677 trast, a small amount of metadata listing or reliably linking the parties, transactions, systems
 678 and provenance of the data would suffice. This relevant context together with the data forms
 679 the basis for accountable analysis on big personal data.

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

684 Some capabilities will likely be essential to all Big Data systems, such as highly scalable
 685 active storage, standard methods for integration with other Big Data systems, and a processing
 686 architecture enabling high speed statistical analytics. But there are and will continue to emerge
 687 multiple types of Big Data systems. Some functions or controls will likely be important —
 688 or even feasible — only for certain types of future systems. For instance, it is reasonable to
 689 expect some systems will specialize in enormous volumes of entirely non-personal data from
 690 many real-time sources (e.g. for soil science, materials engineering, astronomy) while other Big
 691 Data systems will hinge upon mass quantities of highly sensitive personal information (e.g. for
 692 clinical medicine, education and lifelong learning, social entertainment).

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

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

700 The four contextual facets of people, interactions, technology, and data provide a sound

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

Big Data, by its nature, represents a new set of business, legal, and technical capabilities and requirements. Most of the worlds systems today are not capable of ingesting, storing, using, or dynamically flowing Big Data with other systems. Considering that a) Big Data is of high value immediately and higher value in the short and long terms, and b) the young but competitive marketplace of Big Data system components, platforms, applications, and other solutions is a hotbed of innovation it can be predicted that a transition to Big Data systems will continue. The key observation is that virtually all Big Data systems have yet to be designed, implemented, customized, or deployed. Institutions that are the current early adopters of todays Big Data system will soon replace those systems and the rest of the world will adopt Big Data systems in phases over time. Based upon this observation, **AS: ??????????????**

9.2 Research on Big Data for Design of Institutions

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

With the coming of Big Data, we are going to be operating very much out of our old, familiar ballpark. These data are often indirect and noisy, and so interpretation of the data requires greater care than usual. Even more importantly, a great deal of the data is about

human behavior, and the questions are ones that seek to connect physical conditions to social outcomes. Until we have a solid, well-proven, and quantitative theory of social physics, we will not be able to formulate and test hypotheses in the way we can when we design bridges or develop new drugs.

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

An example of such a Living Lab is the ‘open data city just launched by one author (Pentland) 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

754 (when aggregated) for the improvement of the community through data commons that enable
755 social network incentives. The ability to share data safely should enable better idea flow among
756 individuals, companies, and government, and we want to see if these tools can in fact increase
757 productivity and creative output at the scale of an entire city.

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

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

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

776 10 Conclusions

777 Our societies today face unprecedented challenges. Solving those problems will require access
778 to the personal data, so we can understand how the society works, how we move around, what
779 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 legislation and demand of the users. Only with such orchestration it will be possible to change the current feudal system of the data ownership and finally put the immense quantities of the collected personal data to good use.

References

1. Binding obligations on User-Managed Access (UMA) participants. Technical Specifications draft-maler-oauth-umatrust-01, Kantara Initiative, July 2013.
2. User-Managed Access (UMA) profile of OAuth2.0. Technical Specifications draft-hardjono-oauth-umacore-08, Kantara Initiative, December 2013.
3. Nadav Aharony, Wei Pan, Cory Ip, Inas Khayal, and Alex Pentland. Social fmri: Investigating and shaping social mechanisms in the real world. *Pervasive and Mobile Computing*,

- 805 7(6):643–659, 2011.
- 806 4. Sinan Aral and Dylan Walker. Identifying influential and susceptible members of social
807 networks. *Science*, 337(6092):337–341, 2012.
- 808 5. Simon Barber, Xavier Boyen, Elaine Shi, and Ersin Uzun. Bitter to Better – how to
809 make Bitcoin a better currency. In *Proceedings Financial Cryptography and Data Security
810 Conference (Lecture Notes in Computer Science Volume 7397)*, pages 399–414, April 2012.
- 811 6. Ellen Barry. Protests in moldova explode, with help of twitter. *New York Times*, 8, 2009.
- 812 7. Nick Bilton. Girls around me: An app takes creepy to a new level. *The New York Times*.
- 813 8. Center for Environmental & Human Toxicology University of Florida. Development of
814 Cleanup Target Levels (CTLs) For Chapter 62-777, F.A.C. Technical report, Division of
815 Waste Management Florida Department of Environmental Protection, February 2005.
- 816 9. Paul Lukowicz Bert Arnrich Cornelia Setz Gerhard Troster David Tacconi, Oscar Mayora
817 and Christian Haring. Activity and emotion recognition to support early diagnosis of
818 psychiatric diseases. pages 100–102. IEEE, 2008.
- 819 10. Yves-Alexandre de Montjoye, César A Hidalgo, Michel Verleysen, and Vincent D Blondel.
820 Unique in the crowd: The privacy bounds of human mobility. *Scientific reports*, 3, 2013.
- 821 11. Yves-Alexandre de Montjoye, Samuel S Wang, Alex Pentland, Dinh Tien Tuan Anh, An-
822 witaman Datta, Kevin W Hamlen, Lalana Kagal, Murat Kantarcioglu, Vaibhav Khadilkar,
823 Kerim Yasin Oktay, et al. On the trusted use of large-scale personal data. *IEEE Data
824 Eng. Bull.*, 35(4):5–8, 2012.
- 825 12. Ralph A. DeMeo and Sarah Meyer Doar. Restrictive covenants as institutional controls
826 for remediated sites: Worth the effort? *The Florida Bar Journal*, 85(2), 2011.
- 827 13. Nathan Eagle and Alex Pentland. Reality mining: sensing complex social systems. *Per-
828 sonal and ubiquitous computing*, 10(4):255–268, 2006.

- 829 14. Jonathan Woetzel et al. Preparing for china's urban billion. 2009.
- 830 15. Florida Department of Environmental Protection - Division of Waste Management. Insti-
831 tutional Controls Procedures Guidance. http://www.dep.state.fl.us/waste/quick_topics/publications/wc/csf/icpg.pdf, June 2012.
832
- 833 16. Kate Greene. Reality mining. *Technology Review*, 2008.
- 834 17. Lev Grossman. Iran protests: Twitter, the medium of the movement. *Time Magazine*,
835 17, 2009.
- 836 18. Aniko Hannak, Piotr Sapiezynski, Arash Molavi Kakhki, Balachander Krishnamurthy,
837 David Lazer, Alan Mislove, and Christo Wilson. Measuring personalization of web search.
838 In *Proceedings of the 22nd international conference on World Wide Web*, pages 527–538.
839 International World Wide Web Conferences Steering Committee, 2013.
- 840 19. Thomas Hardjono, Patrick Deegan, and John Clippinger. On the Design of Trustworthy
841 Compute Frameworks for Self-Organizing Digital Institutions. In *Proceedings of the 16th*
842 *International Conference on Human-Computer Interaction*, 2014.
- 843 20. Thomas Hardjono, Daniel Greenwood, and Alex Pentland. Towards a trustworthy digital
844 infrastructure for core identities and personal data stores. In *Proceedings of the ID360*
845 *Conference on Identity*. University of Texas, April 2013.
- 846 21. Juniper Networks. Secure Data Access Anywhere and Anytime: Current Landscape and
847 Future Outlook of Enterprise Mobile Security. A forrester consulting thought leadership
848 paper commissioned by att and juniper networks, Forrester Research, October 2012.
- 849 22. Meglena Kuneva. Roundtable on Online Data Collection, Targeting and Profiling . http://europa.eu/rapid/press-release_SPEECH-09-156_en.htm, 2009.
850
- 851 23. David Lazer, Alex Sandy Pentland, Lada Adamic, Sinan Aral, Albert Laszlo Barabasi,
852 Devon Brewer, Nicholas Christakis, Noshir Contractor, James Fowler, Myron Gutmann,

- et al. Life in the network: the coming age of computational social science. *Science (New York, NY)*, 323(5915):721, 2009.
24. Antonio Lima, Manlio De Domenico, Veljko Pejovic, and Mirco Musolesi. Exploiting cellular data for disease containment and information campaigns strategies in country-wide epidemics. School of computer science university of birmingham technical report csr-13-01, University of Birmingham, May 2013.
25. Anmol Madan, Manuel Cebrian, David Lazer, and Alex Pentland. Social sensing for epidemiological behavior change. In *Proceedings of the 12th ACM international conference on Ubiquitous computing*, pages 291–300. ACM, 2010.
26. AC Madrigal. Dark social: We have the whole history of the web wrong. *The Atlantic*, 2013.
27. Alan Mislove, Sune Lehmann, Yong-Yeol Ahn, Jukka-Pekka Onnela, and J Niels Rosenquist. Pulse of the nation: Us mood throughout the day inferred from twitter. *Accessed November, 22(2011):2011*, 2010.
28. Arvind Narayanan and Vitaly Shmatikov. Robust de-anonymization of large sparse datasets. In *Security and Privacy, 2008. SP 2008. IEEE Symposium on*, pages 111–125. IEEE, 2008.
29. Wei Pan, Yaniv Altshuler, and Alex Sandy Pentland. Decoding social influence and the wisdom of the crowd in financial trading network. In *Privacy, Security, Risk and Trust (PASSAT), 2012 International Conference on and 2012 International Conferenece on Social Computing (SocialCom)*, pages 203–209. IEEE, 2012.
30. Wei Pan, Gourab Ghoshal, Coco Krumme, Manuel Cebrian, and Alex Pentland. Urban characteristics attributable to density-driven tie formation. *Nature communications*, 4, 2013.

- 877 31. ALEX PENTLAND. Reality mining of mobile communications: Toward a new deal on
878 data. *The Global Information Technology Report 2008–2009*, page 1981, 2009.
- 879 32. Alex Pentland, David Lazer, Devon Brewer, and Tracy Heibeck. Using reality mining to
880 improve public health and medicine. *Stud Health Technol Inform*, 149:93–102, 2009.
- 881 33. Vivek K Singh, Laura Freeman, Bruno Lepri, and Alex Sandy Pentland. Classifying
882 spending behavior using socio-mobile data. *HUMAN*, 2(2):pp–99, 2013.
- 883 34. Chaoming Song, Zehui Qu, Nicholas Blumm, and Albert-László Barabási. Limits of
884 predictability in human mobility. *Science*, 327(5968):1018–1021, 2010.
- 885 35. Stan Stalnaker. The Ven currency, 2013. <http://www.ven.vc>.
- 886 36. Latanya Sweeney. Simple demographics often identify people uniquely. *Health (San Fran-*
887 *cisco)*, pages 1–34, 2000.
- 888 37. The White House. National Strategy for Trusted Identities in Cyberspace: Enhancing On-
889 line Choice, Efficiency, Security, and Privacy. The White House, April 2011. Available on
890 http://www.whitehouse.gov/sites/default/files/rss_viewer/NSTICstrategy_041511.pdf.
- 891 38. United States Environmental Protection Agency. Institutional Controls Bibliography.
892 <http://www.epa.gov/superfund/policy/ic/guide/biblio.pdf>, December 2005.
- 893 39. United States Environmental Protection Agency. RCRA Corrective Action Institu-
894 tional Controls - glossary. [http://www.epa.gov/epawaste/hazard/correctiveaction/](http://www.epa.gov/epawaste/hazard/correctiveaction/resources/guidance/ics/glossary1.pdf)
895 [resources/guidance/ics/glossary1.pdf](http://www.epa.gov/epawaste/hazard/correctiveaction/resources/guidance/ics/glossary1.pdf), 2007.
- 896 40. United States Environmental Protection Agency. Institutional Controls: A Guide to Plan-
897 ning, Implementing, Maintaining, and Enforcing Institutional Controls at Contaminated
898 Sites. Technical Report OSWER 9355.0-89 EPA-540-R-09-001, EPA, December 2012.

- 899 41. Jessica Vitak, Paul Zube, Andrew Smock, Caleb T Carr, Nicole Ellison, and Cliff Lampe.
900 It's complicated: Facebook users' political participation in the 2008 election. *CyberPsy-*
901 *chology, behavior, and social networking*, 14(3):107–114, 2011.
- 902 42. World Economic Forum. Personal Data: The Emergence of a New
903 Asset Class, 2011. Available on [http://www.weforum.org/reports/](http://www.weforum.org/reports/personal-data-emergence-new-asset-class)
904 [personal-data-emergence-new-asset-class](http://www.weforum.org/reports/personal-data-emergence-new-asset-class).