

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

3 Daniel "Dazza" Greenwood<sup>1,\*</sup>, Arkadiusz Stopczynski<sup>1,2</sup>, Brian Sweatt<sup>1</sup>, Thomas Hardjono<sup>1</sup>,  
4 Alex Sandy Pentland<sup>1</sup>

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)**

10 To realize the promise and prospects of a Big Data society and avoid its security and confiden-  
11 tiality perils, institutions are updating operational frameworks governing business, legal, and  
12 technical dimensions of their internal organization and interactions with the outside world. In  
13 this chapter we explore the emergence of the Big Data Society, outline ways to support it in the  
14 context of institutional controls, and describe future directions of research and development.

15 The control points traditionally relied upon as part of corporate governance, management  
16 oversight, legal compliance, and enterprise architecture must evolve and expand to match op-  
17 erational frameworks for Big Data. An operational framework used for a Big Data-driven or-  
18 ganization requires a balanced set of institutional controls. These institutional controls must  
19 support and reflect greater user control over personal data, and large scale interoperability for  
20 data sharing between and among institutions. Core capabilities of these controls include re-  
21 sponsive rule-based systems governance and fine-grained authorizations for distributed rights  
22 management.

23 Sustaining a healthy, safe, and efficient society is a scientific and engineering challenge go-  
24 ing back to the 1800s, when the Industrial Revolution spurred rapid urban growth, creating  
25 huge social and environmental problems. The remedy then was to build centralized networks

26 that delivered clean water and safe food, enabled commerce, removed waste, provided energy,  
27 facilitated transportation, and offered access to centralized healthcare, police, and educational  
28 services. Those networks formed the backbone of the society as we know it today.

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

34 It does not have to be this way. We can have cities that are protected from pandemics, energy  
35 efficient, have secure food and water supplies, and have much better government. To reach these  
36 goals, however, we need to radically rethink our approach. Rather than static fixed systems,  
37 separated by function — water, food, waste, transport, education, energy — we must consider  
38 them as dynamic, data-driven networks. Instead of focusing only on access and distribution,  
39 we need the networked and self-regulating systems, driven by the needs and preferences of the  
40 citizens. Finally, we need to create the channels for the society to agree upon and communicate  
41 those needs.

42 To ensure a sustainable future society, we must use our new technologies to create a *nervous*  
43 *system* maintaining the stability of government, energy, and public health systems around the  
44 globe. Our digital feedback technologies are today capable of creating a level of dynamic re-  
45 sponsiveness that our larger, more complicated modern society requires. We must reinvent the  
46 systems of the societies within a control framework: sensing the situation, combining these obser-  
47 vations with models of demand and dynamic reaction, and finally using the resulting predictions  
48 to tune the system to match the demands.

49 The engine driving this nervous system is Big Data: the newly ubiquitous digital data, now  
50 available about all aspects of human life. We can analyze patterns of human experience and  
51 ideas exchange within the *digital breadcrumbs* that we all leave behind as we move through the  
52 world: call records, credit card transactions, GPS location fixes, among others. By recording

our choices, these data tell the story of our lives. And this may be very different from what we decide to put on Facebook or Twitter; our postings there are what we choose to tell people, edited according to the standards of the day and filtered to match the persona we are building. Mining social networks can give some great insights about human nature [?, ?, ?]; who we really are is however even more accurately determined by where we spend our time and which things we buy, rather than just what we say we do [?].

The process of analyzing the patterns within these digital breadcrumbs is called reality mining [?, ?], and through it we can learn an enormous amount about who we are. The Human Dynamics research group at MIT have found that we can use them to tell if we are likely to get diabetes [?], or whether we are the sort of person who will pay back loans [?]. 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 [?]. For this reason the magazine Technology Review named our development of reality mining as one of the ten technologies that will change the world [?].

## 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” [?]. 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

79 insurance rates based on their shopping history, or create problems for the entire society in  
80 longer run, for example by limiting user choices and closing them into information bubbles [?].  
81 To achieve the positive possibilities of the new society, we require the *New Deal on Data*, workable  
82 guarantees that the data needed for public good are readily available while at the same time  
83 protecting the citizenry [?].

84 The key insight that motivates the idea of the New Deal on Data is that our data are worth  
85 more when shared, because these aggregated data — averaged, combined across population, and  
86 often distilled to high-level features — inform improvements in systems such as public health,  
87 transportation, and government. For instance, we have demonstrated that data about the way  
88 we behave and where we go can be used to minimize the spread of infectious disease [?,?]. Our  
89 research has reported how we were able to use these digital breadcrumbs to track the spread of  
90 influenza from person to person on an individual level. And if we can see it, we can stop it.

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

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

104 Unfortunately, today most personal data are siloed off in private companies and therefore  
105 largely unavailable. Private organizations collect the vast majority of the personal data in the

106 form of mobility patterns, financial transactions, phone and Internet communications. These  
 107 data must not remain the exclusive domain of private companies, because then they are less  
 108 likely to contribute to the common good. These private organizations must be thus the key  
 109 players in the New Deal on Data framework for privacy and data control. Likewise, these data  
 110 should not become the exclusive domain of the government, as this will not serve the public  
 111 interest of transparency; we should be suspicious of trusting the government with such power.  
 112 Ultimately, the entities who should be empowered to share and make decisions about their data,  
 113 are people themselves: users, participants, citizens.

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

### 125 **3 Personal Data: Emergence of a New Asset Class (Thomas)**

126 It has long been recognized that the first step to promoting liquidity in land and commodity  
 127 markets is to guarantee ownership rights so that people can safely buy and sell. Similarly, the  
 128 first step toward creating greater idea and idea flow ('idea liquidity') is to define ownership rights.  
 129 The only politically viable course is to give individual citizens rights over data that are about  
 130 them and in fact, in the European Union these rights flow directly from the constitution **AS:**  
 131 **Citation? There is no 'EU constitution' per se.** . We need to recognize personal data

132 as a valuable asset of the individual that is given to companies and government in return for  
133 services.

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

- 136 • You have the right to possess data about you. Regardless of what entity collects the data,  
137 the data belong to you, and you can access your data at any time. Data collectors thus  
138 play a role akin to a bank, managing the data on behalf of their customers.
- 139 • You have the right to full control over the use of your data. The terms of use must be opt-  
140 in and clearly explained in plain language. If you are not happy with the way a company  
141 uses your data, you can remove the data, just as you would close your account with a bank  
142 that is not providing satisfactory service.
- 143 • You have the right to dispose of or distribute your data. You have the option to have data  
144 about you destroyed or redeployed elsewhere.

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

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

153 In 2007, one author (Pentland) first proposed the New Deal on Data to the World Economic  
154 Forum [?]. Since then, this idea has run through various discussions and eventually helped shape  
155 the 2012 Consumer Data Bill of Rights in the United States, along with a matching declaration  
156 on Personal Data Rights in the EU. These new regulations hope to accomplish the combined  
157 trick of breaking data out of the current silos, thus enabling public goods, while at the same

time giving individuals greater control over data about them. But, of course this is still a work in progress and the battle for individual control of personal data rages onward.

The World Economic Forum (WEF) has dubbed personal data as the “New Oil” or resource of the 21st century [?]. The discovery of oil and the subsequent development of the oil industry over the past 100 years has spurred not only the development of the automobile industry but also the creation of the global transportation infrastructure, including the massive freeway networks that we see today in the developed nations. The “personal data sector” of the economy today is still in its infancy, its state akin to the oil industry at the late 1890s prior to the development of the Model-T Ford automobile. The productive collaboration between the Government (building the state owned freeways), the private sector (mining and refining oil, building automobiles) and the citizen (the user-base of these services) allowed the developed nations to expand its economies by creating new markets adjacent to the automobile and oil industries.

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

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

For many participants, the risks and liabilities exceed the economic returns. Besides not having the infrastructure and tools to manage personal data, many end-users simply do not see

the benefit of fully participating in the ecosystem. The current focus of many Internet-based service providers is to capture as much personal data from the end-user and to sell this data into the advertising industry. Personal privacy concerns are thus inadequately addressed at best, or simply overlooked in the majority of the cases. The current technologies and laws fall short of providing the legal and technical infrastructure needed to support a well-functioning digital economy.

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

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

- Alignment of key stakeholders: Citizens, the private sector and the public sector need to work in support of one another. Efforts such as NSTIC [?] — albeit still in its infancy —



212 represents a promising direction for a global collaboration.

- 213 • Viewing “data as money”: There needs to be a new change in mindset where an individual’s  
 214 personal data items are viewed and treated in the same way as their money. These personal  
 215 data items would reside in an “account” (like a bank account) where it would be controlled,  
 216 managed, exchanged and accounted for just like personal banking services operate today.
- 217 • End-user centricity: All entities in the ecosystem need to recognize that end-users are  
 218 vital and independent stakeholders in the co-creation and value exchange of services and  
 219 experiences. Efforts such as the *User managed Access* (UMA) initiative [?] point in the  
 220 right direction by designing systems that are user-centric and managed by the user.

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

## 228 4 Enforcing the New Deal on Data (Dazza)

229 How can we enforce this New Deal? The threat of legal action alone is important, but insufficient,  
 230 because if you cannot see abuses then you cannot prosecute them. Moreover, who wants more  
 231 lawsuits anyway? Enforcement can be addressed in significant ways without prosecution of  
 232 public statute or regulation at all. In many fields, companies and governments rely upon multi-  
 233 party frameworks of agreed rules governing common business, legal, and technical practices to  
 234 create effective self-organization and enforcement. These approaches hold promise as a method  
 235 for using institutional controls to form a reliable operational framework balancing the needs for  
 236 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

264 multiparty shared architectural and contractual rules that share the generic characteristic of cre-  
265 ating binding obligations and enforceable expectations on all participants in scalable networks.  
266 Another common characteristic of the system rules design pattern is that the participants in  
267 the network can be widely distributed across very heterogeneous business ownership boundaries,  
268 legal governance structures, and technical security domains. Yet, the parties need not agree to  
269 conform all or most aspects of their basic roles, relationships, and activities in order to connect  
270 to to systems of a trust network. Cross-domain trusted systems must, by their nature, focus  
271 mandatory and enforceable rules narrowly upon the critical items that must be commonly agreed  
272 in order for that network to achieve it's purpose.

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

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

289 The openPDS is a consumer version of a personal cloud trust network that we are now  
290 testing with a variety of industry and government partners. Soon, sharing your personal data

291 could become as safe and secure as transferring money between banks.

292 The Human Dynamics Lab has applied the system rules approach to development of in-  
 293 tegrated business, technical architecture, and rules large scale institutional use of personal  
 294 data stores, available as an example under MIT’s creative commons license by MIT, at [https:](https://github.com/HumanDynamics/SystemRules)  
 295 [//github.com/HumanDynamics/SystemRules](https://github.com/HumanDynamics/SystemRules).

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

## 302 5 Transitioning End-User Assent Practices (Arek)

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

312 Increasing the amount of data controlled by the user and granularity of this control is mean-  
 313 ingless if it cannot be exercised in an informed way. For many years, the End User License  
 314 Agreements (EULAs), long incomprehensible texts have been accepted blindly by the user,  
 315 trusting they have not agreed to anything that could harm them. The process of granting the  
 316 authorizations cannot be too complex, as it would prevent the user from understanding her deci-

sions. At the same time, it cannot be too simplistic, as it may not sufficiently convey the weight of the privacy-related decisions. It is a challenge in itself, to build the end-user assent systems that allow the user to understand and adjust their privacy settings. Complex EULAs do not promote the privacy of the users, effectively pushing them to press *I Agree* in every presented window.

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

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

## 343 6 Business, Legal, and Technical Dimensions of Big Data Sys- 344 tems (Dazza)

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

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

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

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

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 (Thomas)

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 [?]:

“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 controls: governmental controls; proprietary controls; enforcement tools; and infor-

394 mational devices.”

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

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

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

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

418 Consider again the analogy of institutional controls in the context of environmental law, and  
419 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 [?] 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) [?] 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-*

447 trols [?]. Also see: *Institutional Controls Bibliography: Institutional Control, Remedy Selection,*  
 448 *and Post-Construction Completion Guidance and Policy, December 2005* [?].

449 When institutional controls would apply to “separately owned neighboring properties” a  
 450 number of issues arise. Engagement with affected third parties, requiring the party responsible  
 451 for site cleanup to use “best efforts” to attain agreement by third parties to institute the relevant  
 452 institutional controls, use of third party neutrals to resolve disagreements regarding the appli-  
 453 cation with institutional control,s or forcing an acquisition of the neighboring land by forcing  
 454 the party responsible to purchase the property of by purchase of the property directly by the  
 455 EPA [?].

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

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

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

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

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

## 499 8 Scenarios of Use in Context (Dazza)

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

### 507 8.1 Example Scenario: Research Systems

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

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

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

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

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

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

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

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

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

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

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

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

## 582 8.2 Scenarios of Use Today, Tomorrow and the Day After

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

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

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

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

618 By contrast, by adopting the New Deal on Data now it is possible to set conditions that  
619 promote prosperity and invention even before the New Deal on Data frameworks are formally  
620 launched. This is because the uncertainly and confusion about the basic premises and expecta-  
621 tions around personal data and identity will be resolved and so investment and risk taking on  
622 a firm foundation can be unleashed. The value of Big Data can be accessed at less direct cost  
623 and lower risk when uncertainties about privacy liability are addressed and significant the new  
624 value is created by enabling wide scale permission based access to personal data and compu-  
625 tations about such data. Adopting use of personal data services in phases, such one economic  
626 sector, transaction type or data type at a time enables access to the lower costs and new value  
627 in a reasonable manner that allows for time to prepare for and stage each phase of adoption.  
628 By staging and phasing the New Deal on Data typical objections to change based on grounds  
629 of cost, disruption or over regulation can be addressed. Policy incentives can further address  
630 these objections, such as allowing safe harbor protections for conduct of organizations operating  
631 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 [?]. 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.

## 657 9.1 Research on Design and Deployment of Big Data Systems

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

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

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

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

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

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

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

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

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

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

## 713 9.2 Research on Big Data for Design of Institutions

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

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

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

735 An example of such a Living Lab is the ‘open data city just launched by one author (Pentland)  
 736 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

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