



## ANALYSIS OF PRIVACY COMPLIANCE OF BOOTSTRAPPING IN BIG DATA SYSTEMS

**Dr. S. A. Bhura**

Assistant Professor, Dept. of Information Technology, Babasaheb Naik College of Engg. Pusad.  
*sabhura@rediffmail.com*

---

### ABSTRACT—

*Bootstrap is a computer approach to get statistical accuracy. It is applied to a wide variety of statistical procedures like non parametric regressions, classification trees or density estimation. This technique requires fewer assumptions and offers greater accuracy and insight than other standard methods for many problems. With the rapid increase in cloud services collecting and using user data to offer personalized experiences, ensuring that these services comply with their privacy policies has become a business imperative for building user trust. This paper mainly focus on two techniques (a) LEGALEASE and (b) GROK that can be use to maintain priavacy in bootstrapping of big data.*

**Keywords—** *Bootstrap, Legalease, Grok-mapper.*

---

### 1. INTRODUCTION

With the rapid increase in cloud services collecting and using user data to offer personalized experiences, ensuring that these services comply with their privacy policies has become a business imperative for building user trust. Security has been one of the major issues required for the use of Big Data. Let us refer to [1] for an overview and analysis of the top ten Big Data security and privacy challenges. However, most compliance efforts in industry today rely on manual review processes and audits designed to safeguard user data, and therefore are resource intensive and lack coverage. Manual reviews and audits in that firstly legal team craft the privacy policy then experts interprets it after that developer writes code for tha and finally audit team verifies compliance which is very time-consuming, resource-intensive, lack coverage, and, thus, inherently do not scale well in large companies ,indeed, there have been cases where internal processes have not caught policy violations. to avoid all this conflicts this paper describes two techniques to achieve privacy compliance , they are (a)LEGALEASE—a language that allows specification of privacy policies that impose restrictions on how user data is handled; and (b) GROK—a data inventory for Map-Reduce-like big data systems that tracks how user data flows among programs. GROK maps code-level schema elements to datatypes in LEGALEASE.encoding policy in LEGALEASE using the GROK data inventory, we decouple interactions so policy specification,interpretation, product development, and continuous auditing can proceed in parallel.

To contextualize the challenges in performing automated privacy compliance checking in a large company with tens of thousands of employees, it is useful to understand the division of labor and responsibilities in current compliance workflows [2], [3]. Privacy policies are typically crafted by lawyers in a corporate legal team to adhere to all applicable laws and regulations worldwide. Due to the rapid change in product features and internal processes, these policies are necessarily specified using high-level policy concepts that may not cleanly

map to the products that are expected to comply with them. For instance, a policy may refer to “IP Address” which is a high-level policy concept, and the product may have thousands of data stores where data derived from the “IP Address” is stored and several thousand processes that produce and consume this data, all of which have to comply with policy. The task of interpreting the policy as applicable to individual products then falls to the tens of privacy champions embedded in product groups. Privacy champions review product features at various stages of the development process, offering specific requirements to the development teams to ensure compliance with policy. The code produced by the development team is expected to adhere to these requirements. Periodically, the compliance team audits development teams to ensure that the requirements are met .

Let us assume that we are interested in checking compliance for an illustrative policy clause that promises “full IP address will not be used for advertising.”. The privacy champion reviewing the algorithm design for, say online advertisement auctions, may learn in a meeting with the development team that they use the IP address to infer the user’s location, which is used as a bid-modifier in the auction. The privacy champion may point out that this program is not compliant with the above policy and suggest to the development team to truncate the IP address by dropping the last octet to comply with the policy, without significantly degrading the accuracy of the location inference. The development team then modifies the code to truncate the IP address. Periodically, the audit team may ask the development team whether the truncation code is still in place. Later, the advertising abuse detection team may need to use the IP address. This may result in a policy exception, but may come with a different set of restrictions, e.g., "IP address may be used for detecting abuse. In such cases it will not be combined with account information." The entire process is highly manual, with each step sometimes taking weeks to identify the right people to talk to and multiple meetings between different groups that may as well be communicating in different languages.

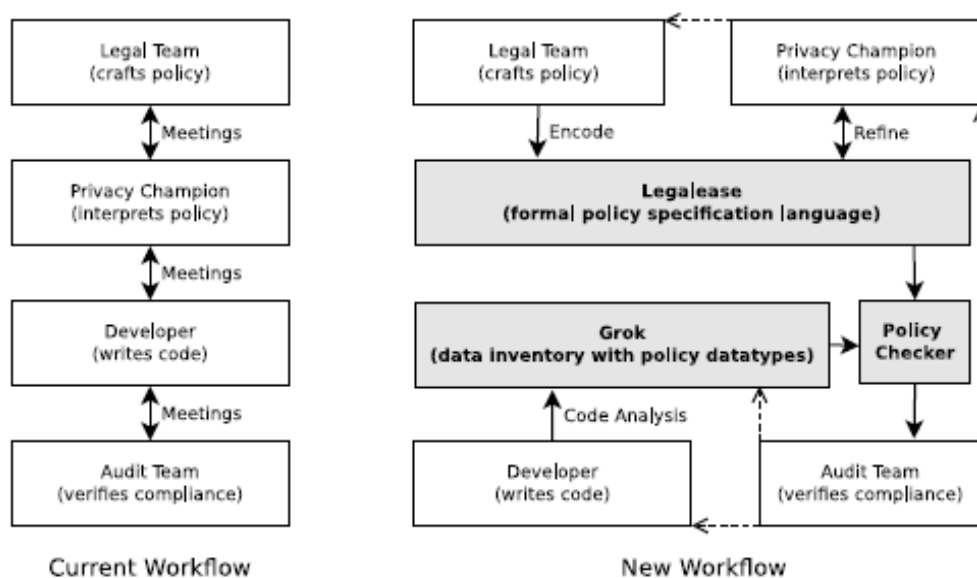


Figure 1: Privacy Compliance Workflow.

The LEGALEASE language. LEGALEASE is an usable, expressive, and enforceable privacy policy language. The primary design criteria for this language were that it (a) be *usable* by the policy authors and privacy

champions; (b) be *expressive* enough to capture real privacy policies of industrial-scale systems, e.g., Bing; (c) and should allow *compositional reasoning* on policies.

## 2. Review of Bootstrap, LEGALEASE and GROK- mapper

### 2.1 Bootstrap

Bootstrapping is an analogy in which the observed data assume the role of an underlying population: ons and confidence intervals are obtained by drawing samples from the empirical sample, as R. Stinewroitein [4]. A typical problem in applied statistics involves the estimation of an unknown parameter. The two main questions are: what estimator should be used? and having chosen a particular one, how accurate is the estimator?

Bootstrap is a general methodology to answer the second question, as stated by Efron and Tibshiraniin [5]. This work is framed within the Intelligent Control Group, Universidad Politécnica de Madrid, whose members are carrying out research into robotics and intelligent control systems. Research covers a wide number of areas: path finding, navigation, speaking, facial expression, mood and knowledge management.

### 2.2 Legalease

As the intended users for LEGALEASE are policy authors and privacy champions with limited training in formal languages, enabling usability is essential. To this end, LEGALEASE enforces syntactic restrictions ensuring that encoded policy clauses are structured very similarly to policy texts. Specifically, building on prior work on a first order privacy logic , policy clauses in LEGALEASE allow (resp. deny) certain types of information flows and are refined through exceptions that deny (resp. allow) some sub-types of the governed information flow types. This structure of nested allow-deny rules appears in many practical privacy policies, including privacy policies for Bing and Google. A distinctive feature of LEGALEASE (and a point of contrast from prior work based on first-order logic and first order-temporal logic ) is that the semantics of policies is compositional: reasoning about a policy is reduced to reasoning about its parts. This form of compositionality is useful because the effect of adding a new clause to a complex policy is locally contained (an exception only refines its immediately enclosing policy clause) To validate the usability of LEGALEASE by its intended users, we conduct a user study among policy writers and privacy champions within Microsoft. On the other hand, by encoding Bing and Google’s privacy policies regarding data usage on their servers, we demonstrate that LEGALEASE retains enough expressiveness to capture real privacy policies of industrialscale systems. We illustrate LEGALEASE with the help of grammer for it that build up to a complex clause. The simplest LEGALEASE policy is DENY. The policy contains a single clause; the clause contains no exceptions and no attribute restrictions. The policy, rather uninterestingly,simply denies everything. We next add a restriction along the DataType attribute for graph nodes to which IP address flows. DENY DataType IPAddress (e.g., that the IP address has been truncated before it can be used).

$$\begin{aligned}
 & \text{Policy Clause } C ::= D \mid A \\
 \text{Deny Clause } D & ::= \text{DENY } T_1 \cdot \cdot T_n \text{ EXCEPT } A_1 \dots A_m \mid \text{DENY } T_1 \dots T_n \\
 \text{Allow Clause } A & ::= \text{ALLOW } T_1 \cdot \cdot T_n \text{ EXCEPT } D_1 \dots D_m \\
 & \quad \mid \text{ALLOW } T_1 \dots T_n \\
 \text{Attribute } T & ::= (\text{attribute-name}) v_1 \cdot \cdot v_l \\
 \text{Value } v & ::= (\text{attribute-value})
 \end{aligned}$$

TABLE 1

## GRAMMER for LEGALEASE

A LEGALEASE policy (Table 1) is rooted in a single toplevel policy clause. A *policy clause* is a layered collection of (alternating) ALLOW and DENY clauses where each clause relaxes or constricts the enclosing clause. Informally, an ALLOW clause permits graph nodes labeled with any subset of the attribute values listed in the clause, and a DENY clause forbids graph nodes labeled with any set that overlaps with the attribute values in the clause. The layering of clauses determines the context within which each clause is checked.

### 2.2.1 Design Goals of LEGALEASE

- **Usability:** Policy clauses in LEGALEASE are structured very similarly to clauses in the English language policy. This correspondence is important because no single individual in a large company is responsible for all policy clauses.
- **Expressivity:** LEGALEASE clauses are built around an attribute abstraction (described below) that allows the language to evolve as policy evolves. For instance, policies today tend to focus on access control, retention times, and segregation of data in storage, [6], [7], [8].
- **Compositional Reasoning :** LEGALEASE provides meaningful syntactic restrictions to allow compositional reasoning where the result of checking the whole policy is a function of reasoning on its parts.

### 2.3 GROK

The GROK mapper. GROK is a data-inventory for MapReduce-like big data systems. It maps every dynamic schemaelement (e.g., members of a tuple passed between mappers and reducers) to datatypes in LEGALEASE. This inventory can be viewed as a mechanism for annotating existing programs written in languages like Hive , Dremel , or Scope with the information flow types (datatypes) in LEGALEASE. Our primary criteria for this technique is (a) be *bootstrapped* with minimal developer effort; (b) reflect *exhaustive and up-to-date* information about all data in the Map-Reduce-like system; and (c) make it easy to *verify* (and update) the mapping from schema-elements to LEGALEASE datatypes. The inventory mappings combine information from a number of different sources each of which has its own characteristic coverage and quality. For instance, syntactic analysis of source code (e.g., applying pattern-matching to column names) has high coverage but low confidence, whereas explicit annotations added by developers has high confidence but low coverage.

## 3.LIMITATIONS

### 3.1 Expressiveness:

LEGALEASE cannot express policies based on first-order temporal-logic. However, LEGALEASE is enough to express privacy policies.

### 3.2 Inference of Sensitive Data:

Unless explicitly labeled, GROK cannot detect inference from non-sensitive data to sensitive data

### 3.3 Precision:

Major source of precision comes from overly conservative treatment of UDF.

### 3.4 False Negatives

The information flow analysis in GROK is conservative , hence bootstrapping that leads to more coverage of the graph with labels.

#### 4.CONCLUSION.

Traditional security methods cannot be applied to big data due to its large volume and variety hence legalease and grok can be more bebefitful for achieving privacy in bootstrapping of big data system. In this paper, we have seen an overview of two techniques for automate privacy policy compliance in bootstrapping of big data concept. Automated privacy compliance checking. LEGALEASE: stating privacy policies as a form of restrictions on information flows. GROK: data inventory that maps low level data types in code to high level policy concepts.

#### 5.REFERENCES

- [1] S. Rajan *et al.* (2013, Oct. 12). *Expanded Top Ten Big Data Security and Privacy Challenges*. Cloud Security Alliance, Los Angeles, CA, USA[Online]. Available: <http://cloudsecurityalliance.org/research/big-data/>
- [2] Shayak Sen, Saikat Guha, Anupam Datta , Sriram K. Rajamani†, Janice Tsai‡ and Jeannette M. Wing “Bootstrapping Privacy Compliance in Big Data” 2014 IEEE Symposium on Security and Privacy .
- [3] International Journal of Artificial Intelligence and Interactive Multimedia, Vol. 1, N° 6.”Improving web learning through model optimization using Bootstrap.
- [4] Stine, R. An Introduction to Bootstrap. Sociological Methods and Research, Vol. 18, Nos. 2&3, November 1989/February 1990 243-291. 1990.ISSN: 1552-8294.
- [5] Effron, B. & Tibshirani, R. J..Bootstrap Methods for Standard Errors, Confidence Intervals, and Other Measures of Statistical Accuracy. Statistical Science, 1986, Vol. 1, 54-77. 1986.ISSN 1726-3328 IEEE Computer Society, 2008, pp. 107–116.
- [6] Facebook. (2012, Dec.) Data use policy. Available: [https://www.facebook.com/full\\_data\\_use\\_policy](https://www.facebook.com/full_data_use_policy)
- [7] Google. (2013, Jun.) Privacy policy. Available: <http://www.google.com/policies/privacy/>
- [8] (2013, Oct.)Bing privacy statement. Microsoft.[Online]. Available:<http://www.microsoft.com/privacystatement/en-gb/bing/default.aspx>