# Bootstrapping Privacy Compliance in Big Data Systems

# Anupam Datta

## Fall 2014

# Privacy Compliance for Bing



## Setting:

▸ Auditor has access to source code

# The Privacy Compliance Challenge

**Legal Team**
Crafts Policy

Meetings

**Privacy Champion**
Interprets Policy

Meetings

**Developer**
Writes Code

Meetings

**Audit Team**
Verifies Compliance

English
Privacy Policy

Compliant?

Millions of Lines of
Undocumented Code

# A Streamlined Audit Workflow

**Legal Team**
Crafts Policy

Encode

**Legal*ease***
A formal policy specification language

Refine

Pr
Interprets Policy

**Grok**
Data inventory with policy labels

Annotated Code

Legalease Policy

D
Wri

Code analysis
Developer annotations

Update Grok

**Checker**

Potential violations

**Audit Team**
Verifies Compliance

Fix code

# A Streamlined Audit Workflow

| | |
|---|---|
| **Legal Team** Crafts Policy | Privacy Champion Interprets Policy |
| **Legal*ease*** A formal policy specification language | |
| **Grok** Data inventory with policy datatypes | Legalease Code Policy Checker |
| **Developer** Writes Code | Audit Team Verifies Compliance |

Encode

Code analysis, developer annotations

Update Grok   Potential violations

Fix code

**Workflow** for privacy compliance

**Legalease,** usable yet formal policy specification language

**Grok,** bootstrapped data inventory for big data systems

**Scalable** implementation for Bing

# A Streamlined Audit Workflow

**Legal Team**
Crafts Policy

**Privacy Champion**
Interprets Policy

Encode

Refine

## Legal*ease*
A formal policy specification language

**Grok**
Data inventory with policy datatypes

Annotated
Code

Legalease
Policy

**Checker**

Code analysis, developer annotations

Update Grok

Potential violations

**Developer**
Writes Code

Fix code

**Audit Team**
Verifies Compliance

# Legalease : Syntax

$$
\begin{array}{lll}
\text{Policy Clause } C & ::= & D \mid A \\
\text{Deny Clause } D & ::= & \text{DENY } T_1 \cdots T_n \text{ EXCEPT } A_1 \cdots A_m \\
& & \mid \text{DENY } T_1 \cdots T_n \\
\text{Allow Clause } A & ::= & \text{ALLOW } T_1 \cdots T_n \text{ EXCEPT } D_1 \cdots D_m \\
& & \mid \text{ALLOW } T_1 \cdots T_n \\
\text{Attribute } T & ::= & \langle \text{attribute-name} \rangle \ v_1 \cdots v_l \\
\text{Value } v & ::= & \langle \text{attribute-value} \rangle
\end{array}
$$

# Legalease

**DENY** *Datatype* IPAddress

　　*UseForPurpose* Advertising

We will **not** use **full IP Address** for **Advertising**.

# Legalease

**DENY** *Datatype* IPAddress
   *UseForPurpose* Advertising
**EXCEPT**
  **ALLOW**
   *Datatype* IPAddress: Truncated
  **ALLOW**
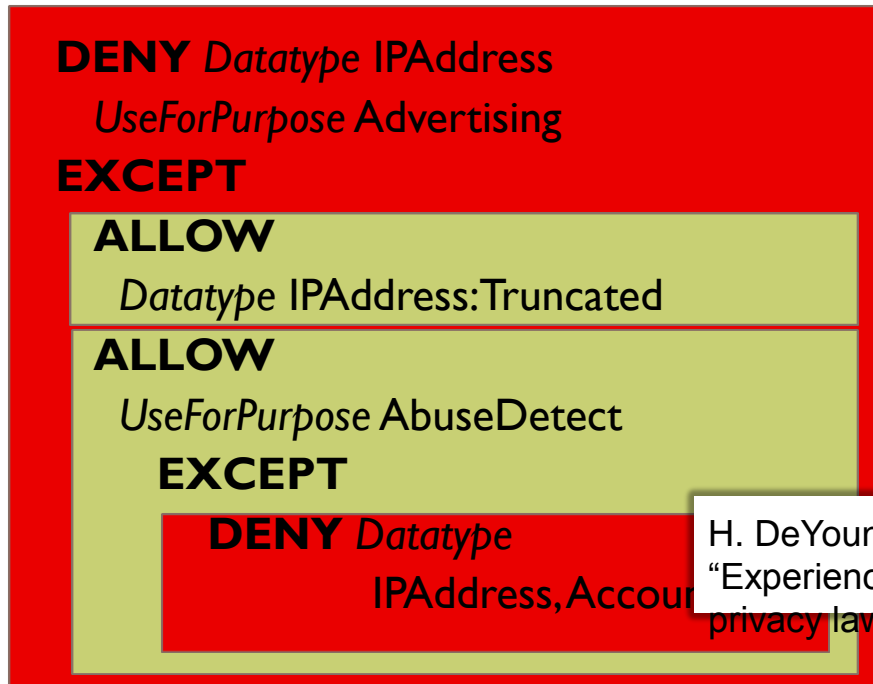   *UseForPurpose* AbuseDetect
    **EXCEPT**
     **DENY** *Datatype*
      IPAddress, AccountInfo

We will **not** use **full IP Address** for **Advertising**. IP Address may be used for **detecting abuse**. In such cases, it will not be combined with **account information.**
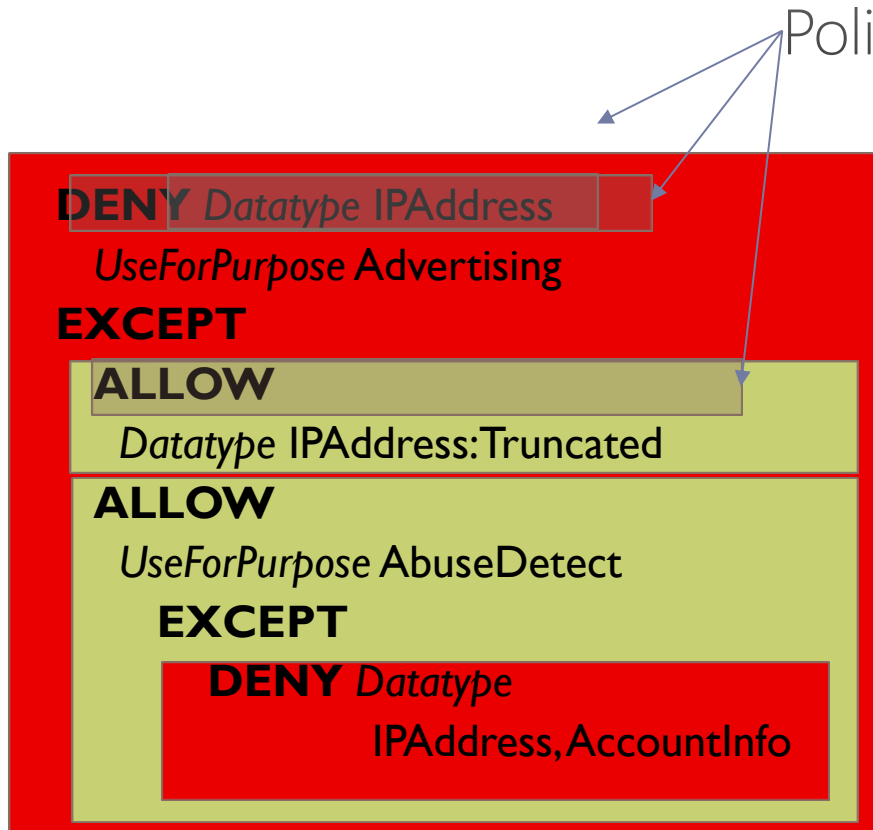
# Designed for Usability

**DENY** *Datatype* IPAddress
   *UseForPurpose* Advertising
**EXCEPT**

> **ALLOW**
>    *Datatype* IPAddress: Truncated
>
> **ALLOW**
>    *UseForPurpose* AbuseDetect
>       **EXCEPT**
>          **DENY** *Datatype*
>                IPAddress, Accou...

## Exceptions
- How legal texts are structured
- One-to one correspondence

## Local Reasoning
- Each exception refines its immediate parent
- Formally proven property

H. DeYoung, D. Garg, L. Jia, D. Kaynar, and A. Datta, "Experiences in the logical specification of the HIPAA and GLBA privacy laws"
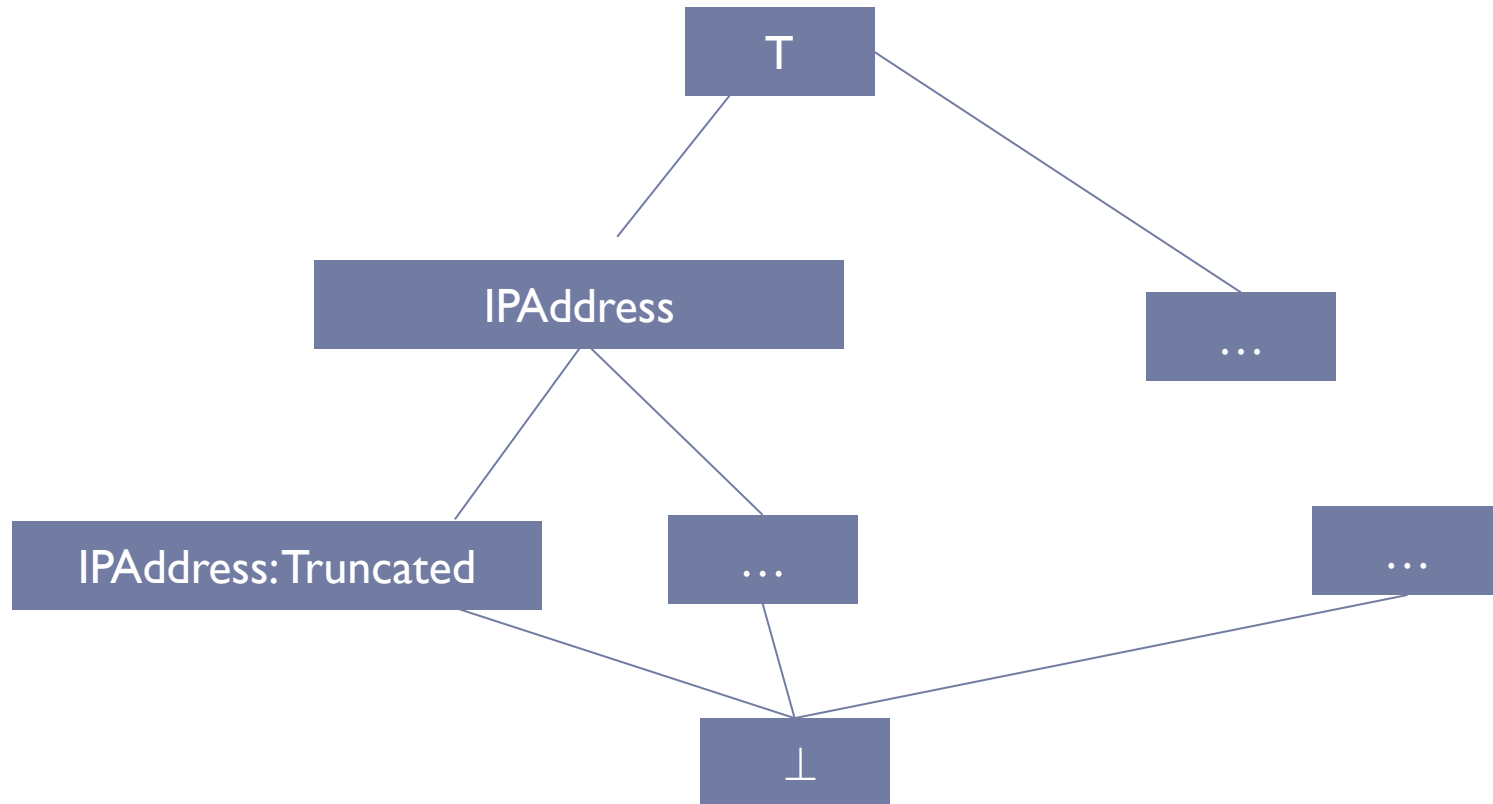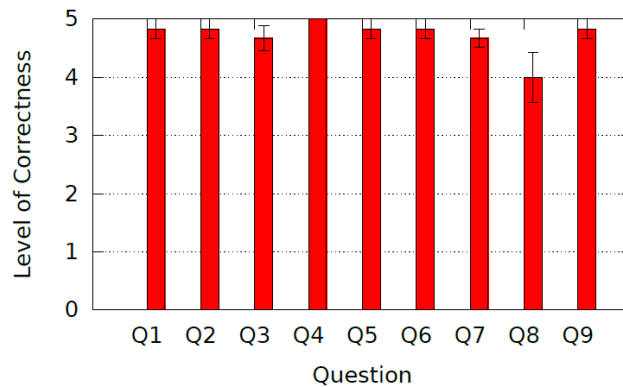
# Legalease : In Action

Policy Labels

Program

*Datatype*: IPAddress, AccountInfo
*UseForPurpose*: AdsAbuseDetection

**DENY** *Datatype* IPAddress
  *UseForPurpose* Advertising
**EXCEPT**
  **ALLOW**
  *Datatype* IPAddress: Truncated
  **ALLOW**
  *UseForPurpose* AbuseDetect
    **EXCEPT**
      **DENY** *Datatype*
          IPAddress, AccountInfo

We will **not** use **full IP Address** for **Advertising**. IP Address may be used for **detecting abuse**. In such cases, it will not be combined with **account information.**
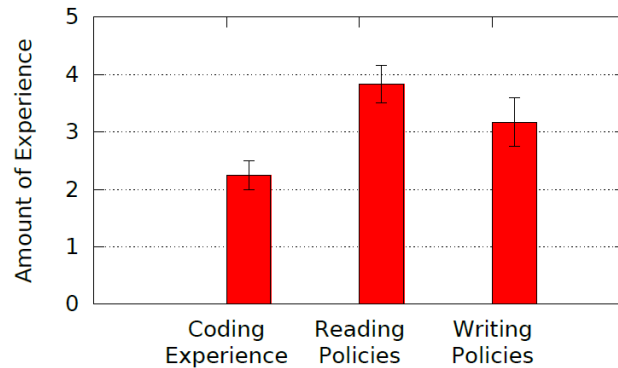
# A Lattice of Policy Labels



- If "IPAddress" use is allowed then so is everything below it
- If "IPAddress:Truncated" use is denied then so is everything above it

# Legalease Usability



Survey taken by 12 policy authors within Microsoft
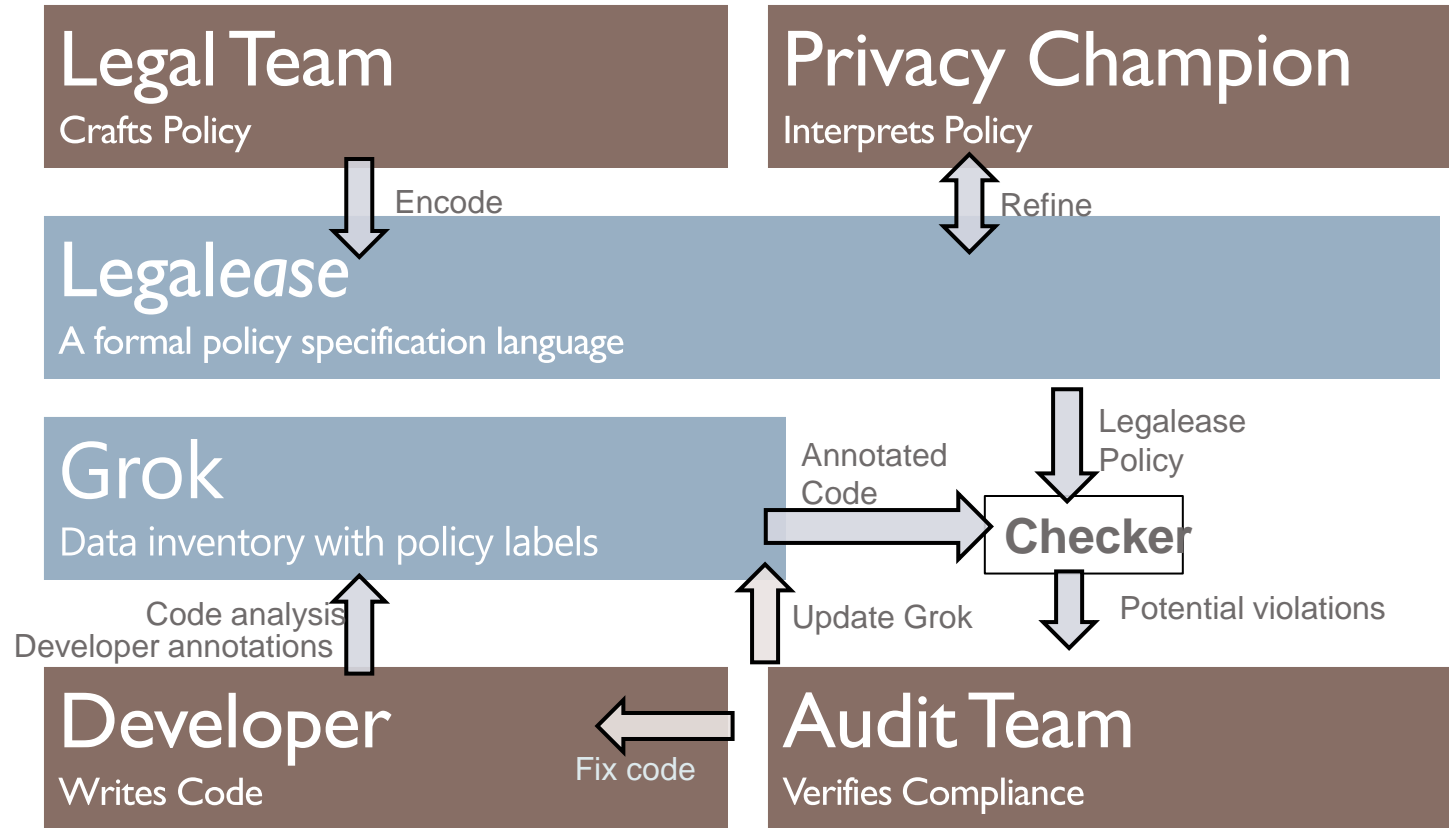
Encode Bing data usage policy after a brief tutorial
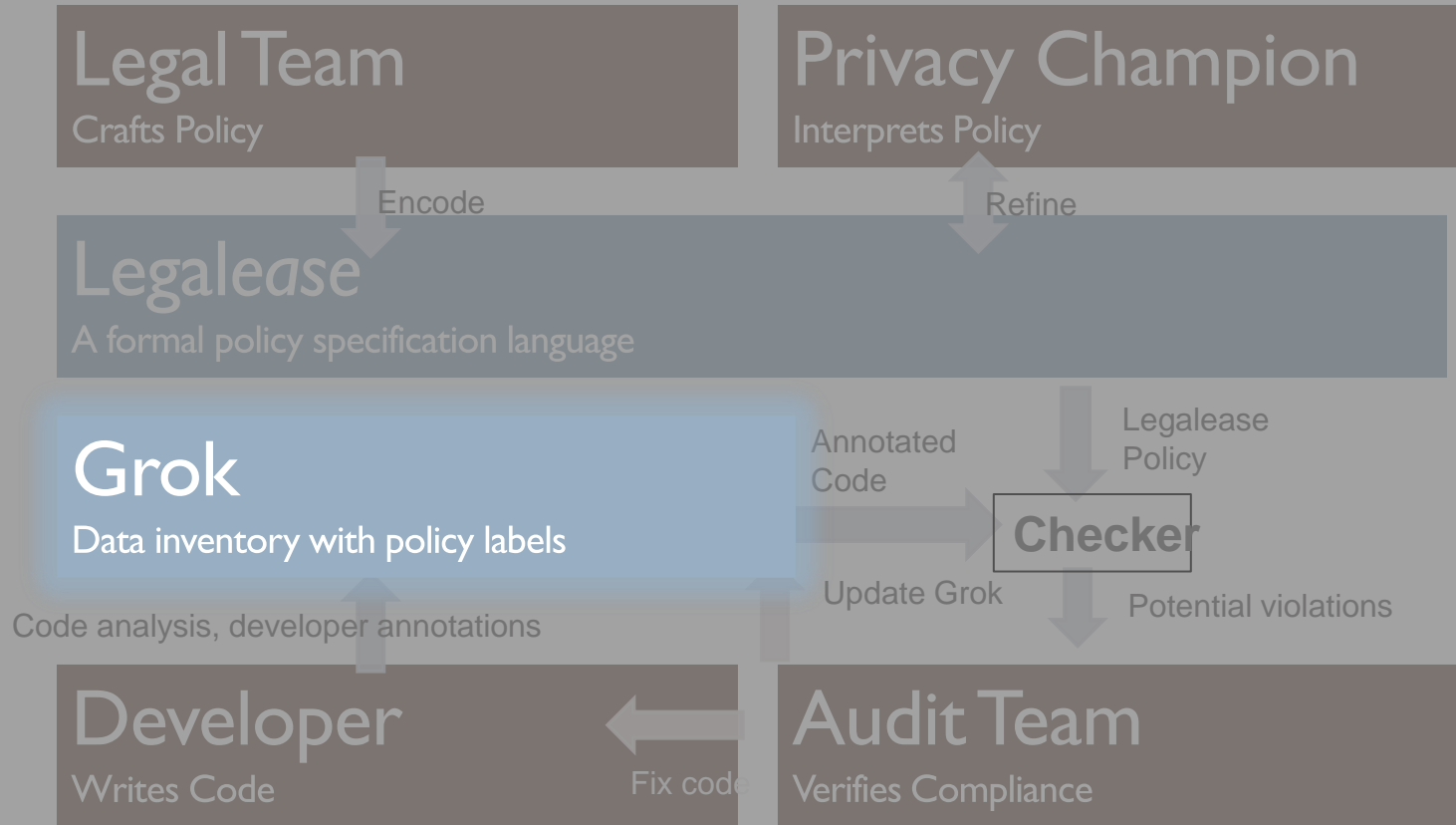
Time spent

2.4 mins on the tutorial

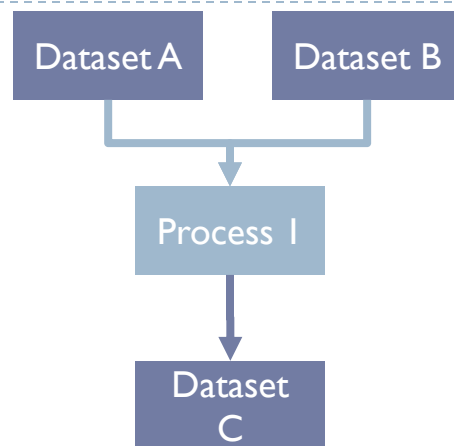14.3 mins on encoding policy

High overall correctness

# A Streamlined Audit Workflow

**Legal Team**
Crafts Policy

**Privacy Champion**
Interprets Policy

Encode

Refine

**Legal*ase***
A formal policy specification language

**Grok**
Data inventory with policy labels

Annotated Code

Legalease Policy

**Checker**

Code analysis
Developer annotations

Update Grok

Potential violations

**Developer**
Writes Code

Fix code

**Audit Team**
Verifies Compliance

# A Streamlined Audit Workflow



**Legal Team**
Crafts Policy

**Privacy Champion**
Interprets Policy

Encode

Refine

**Lega*se***
A formal policy specification language

**Grok**
Data inventory with policy labels

Annotated
Code

Legalease
Policy

**Checker**

Code analysis, developer annotations

Update Grok

Potential violations

**Developer**
Writes Code

**Audit Team**
Verifies Compliance

Fix code

# Map-Reduce Programming Systems

| Dataset A | Dataset B |
|-----------|-----------|

Process 1

Dataset C

Scope, Hive, Dremel

Data in the form of Tables

Code Transforms Columns to Columns
  No Shared State
  Limited Hidden Flows

```
users =
    SELECT _name, _age FROM datasetAB
user_tag =
    SELECT GenerateTag(_name, _age)
        FROM users
OUTPUT user_tag TO datasetC
```
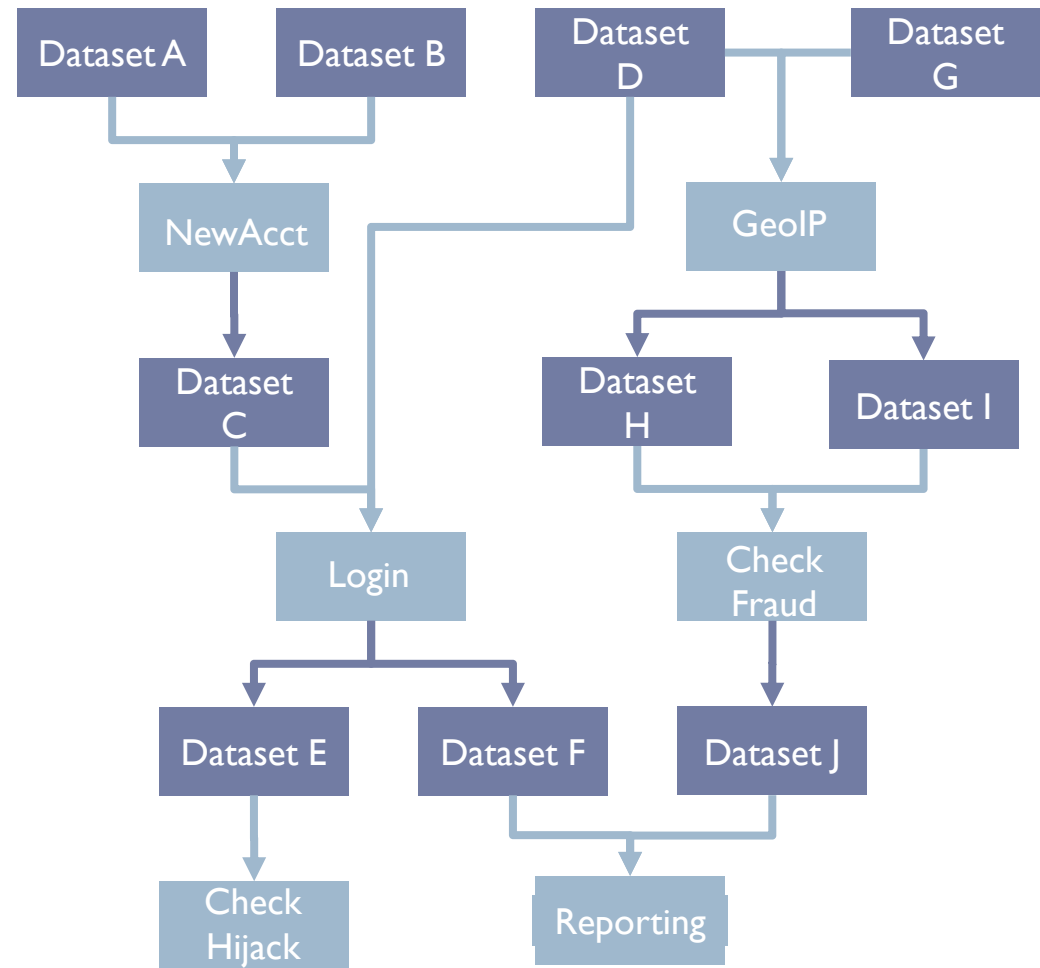
# Grok

# Grok

## Purpose Labels
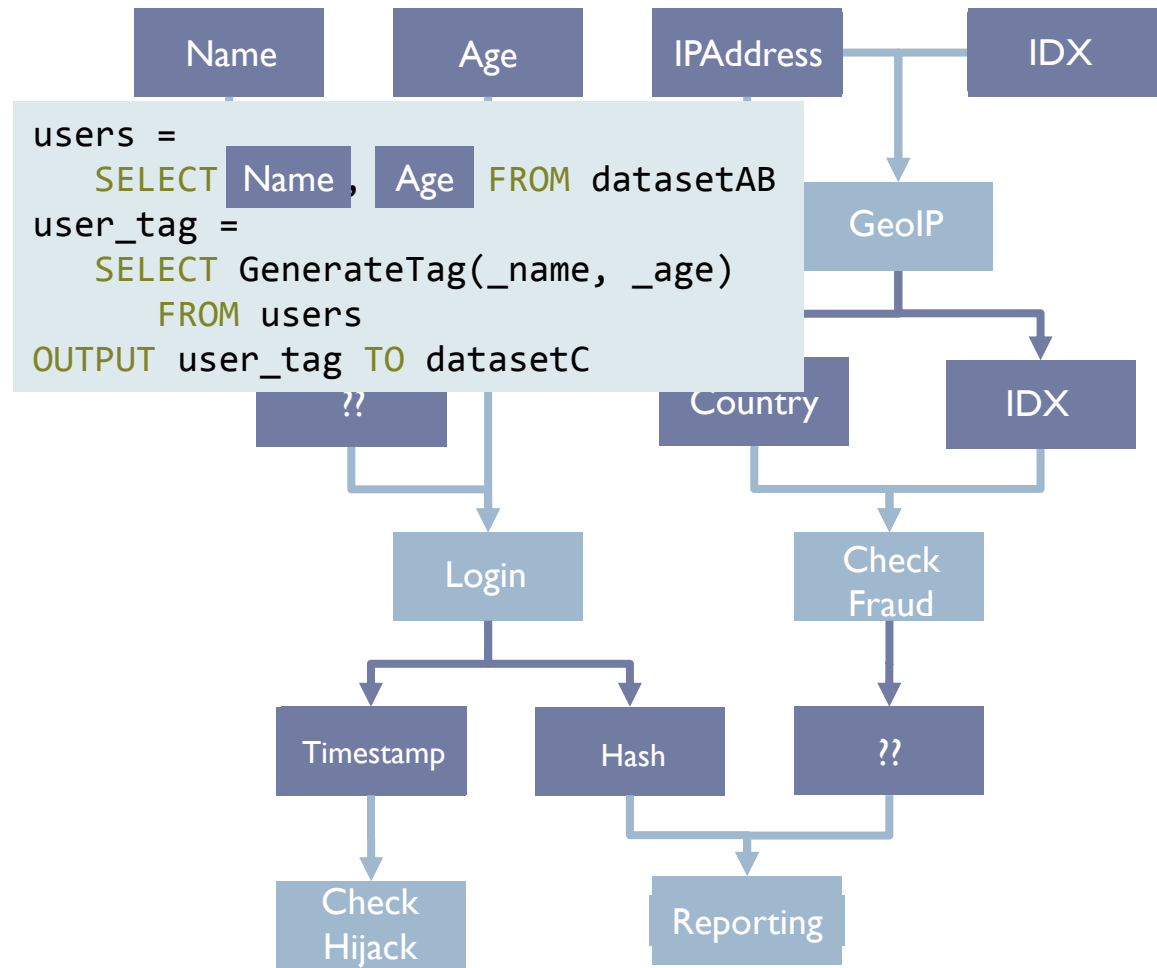
Annotate programs with purpose labels

# Grok



**Purpose Labels**

Annotate programs with purpose labels

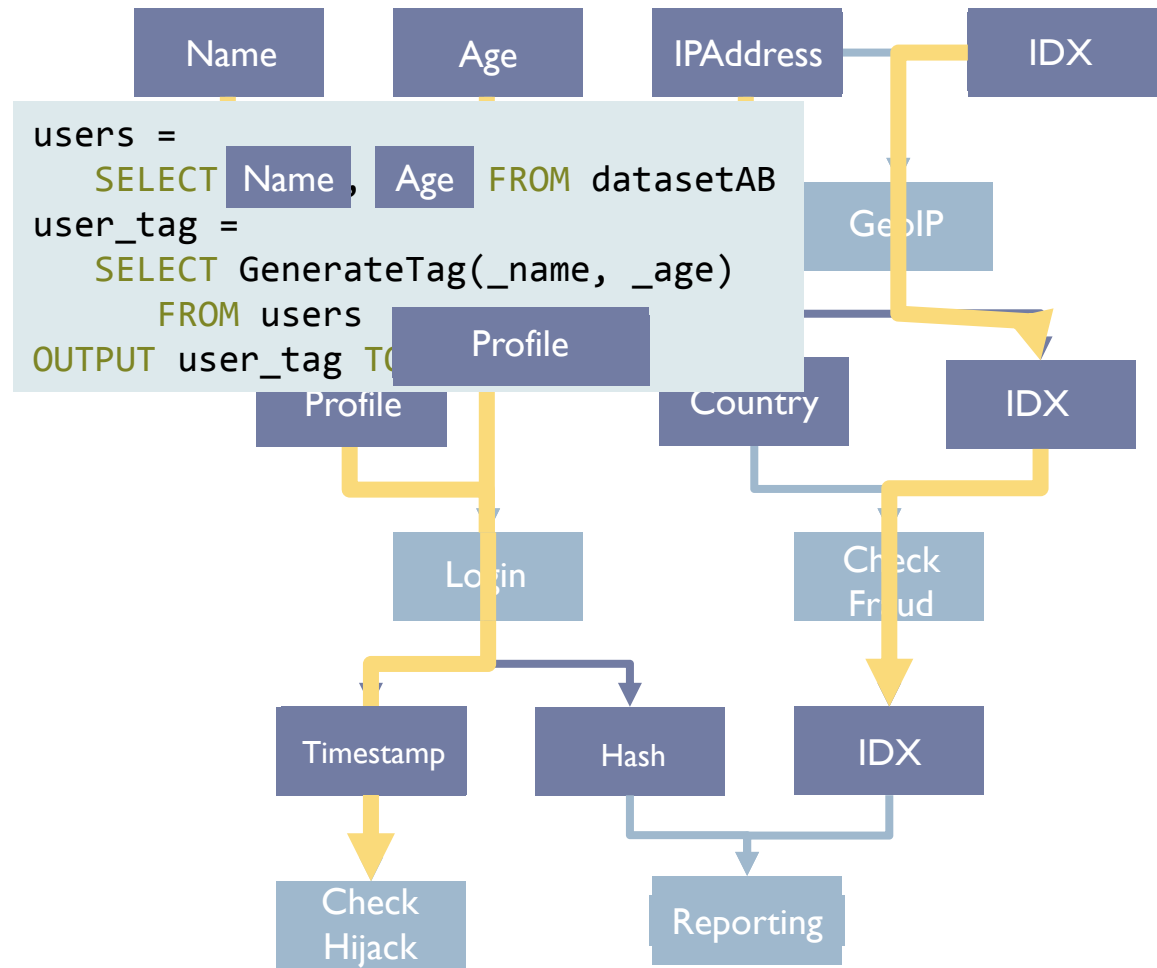**Initial Data Labels**

Heuristics and Annotations

```
users =
    SELECT Name , Age  FROM datasetAB
user_tag =
    SELECT GenerateTag(_name, _age)
        FROM users
OUTPUT user_tag TO datasetC
```

# Grok

**Purpose Labels**

Annotate programs with purpose labels

**Initial Data Labels**

Heuristics and Annotations

**Flow Labels**

Source labels propagated via data flow graph



```
users =
    SELECT Name , Age  FROM datasetAB
user_tag =
    SELECT GenerateTag(_name, _age)
        FROM users
OUTPUT user_tag TO
```

D. E. Denning. "A lattice model of secure information flow"

# A Lattice of Policy Labels



- If "Profile" use is allowed then so is everything below it
- If "Name" use is denied then so is everything above it

# Implicit flows

```
users =
    SELECT Name , Age  FROM datasetAB

users_35 =
    SELECT _name
        FROM users
        WHERE (_age > 35)

OUTPUT users_35 TO    Profile
```

**Beyond direct flows discussed in healthcare audit examples**

# Map-Reduce

## Map

Operate on rows
in parallel
eg. filtering

## Reduce

Combine groups of rows
eg. aggregation

```
users =
    SELECT  Name ,  Age   FROM datasetAB

users_35 =
    SELECT _name, _age
        FROM users
        WHERE (_age > 35)

ages_35 =
    SELECT _age, COUNT(_name) AS  Profile
        FROM users_35
        GROUP BY _age

OUTPUT ages_35 TO datasetC
```

# Combine Noisy Sources

Carefully curated
regular expressions

Leverages developer
conventions

Significant Noise

**Variable Name
Analysis**

Expensive
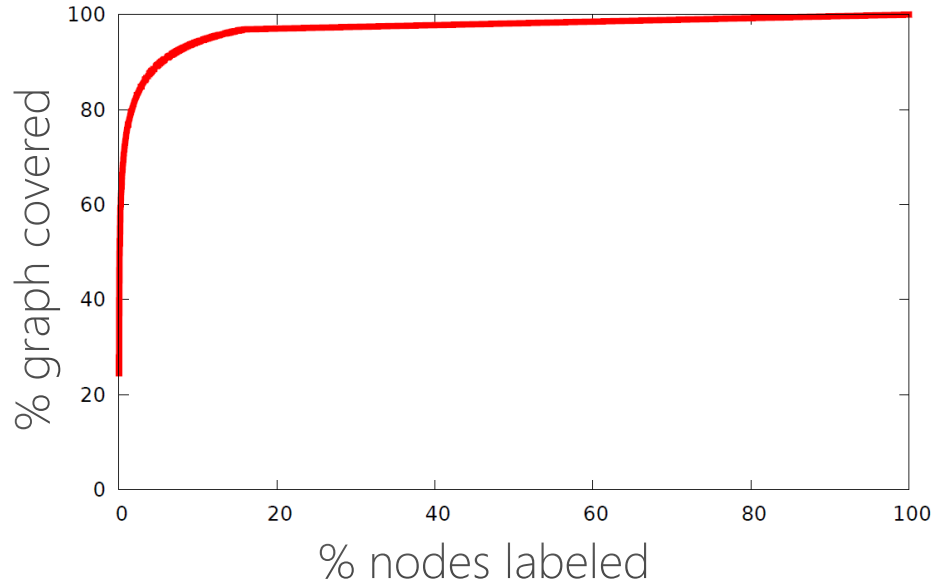
Low Noise

**Developer
Annotations**

Very Expensive

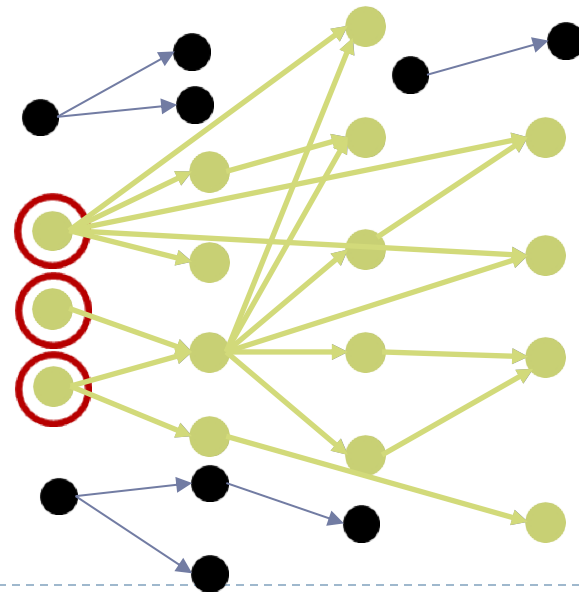Definitive

Need very few of these

**Auditor
Verification**

# Why Bootstrapping Grok Works



% graph covered vs % nodes labeled

Pick the nodes which will label the most of the graph

~200 annotations label 60% of nodes

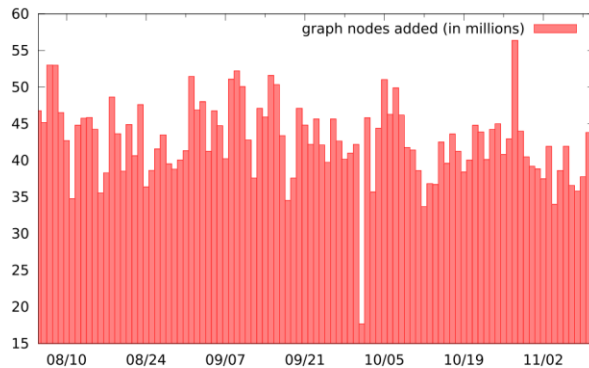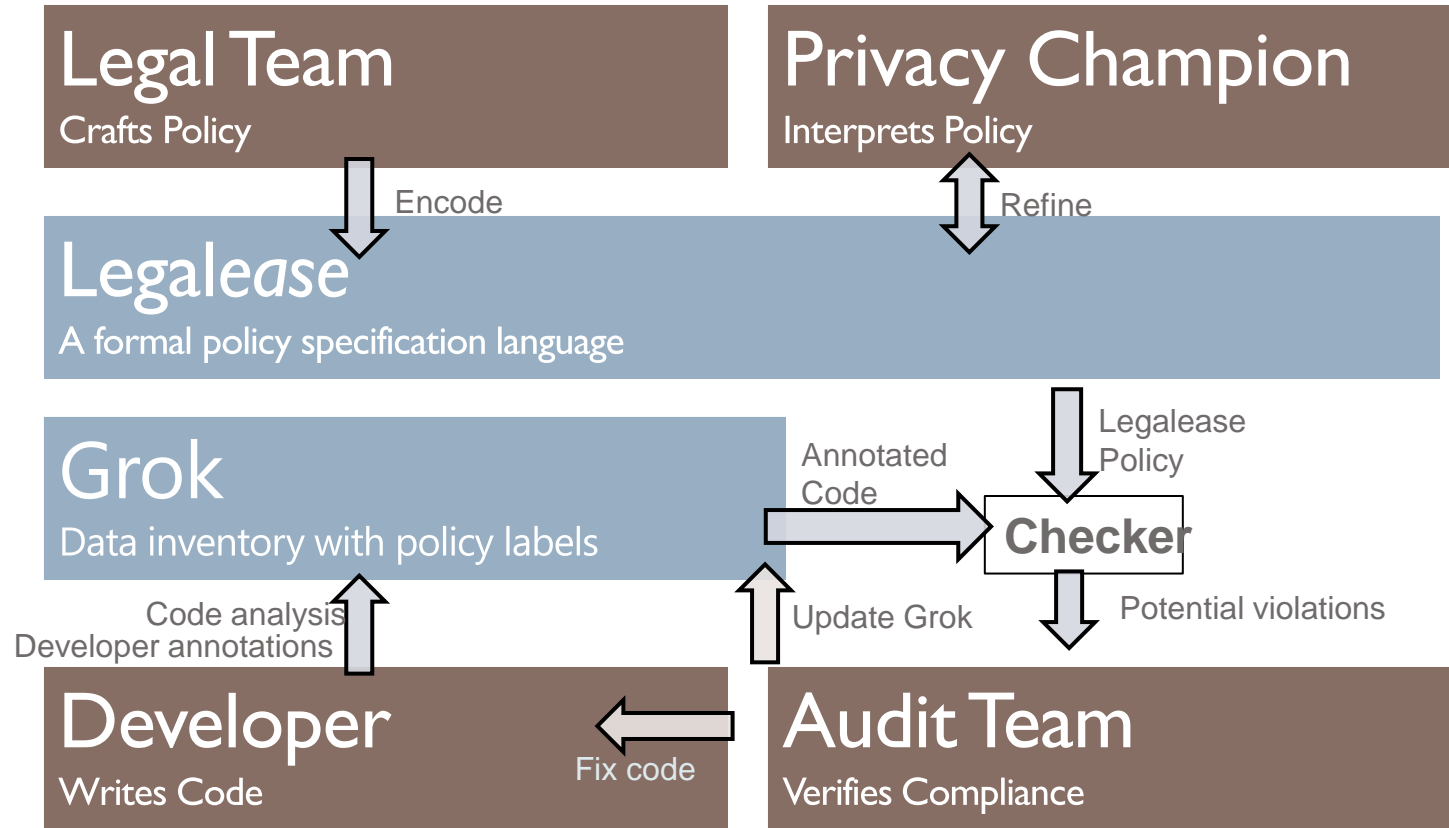A small number of annotations is enough to get off the ground.

# Scale
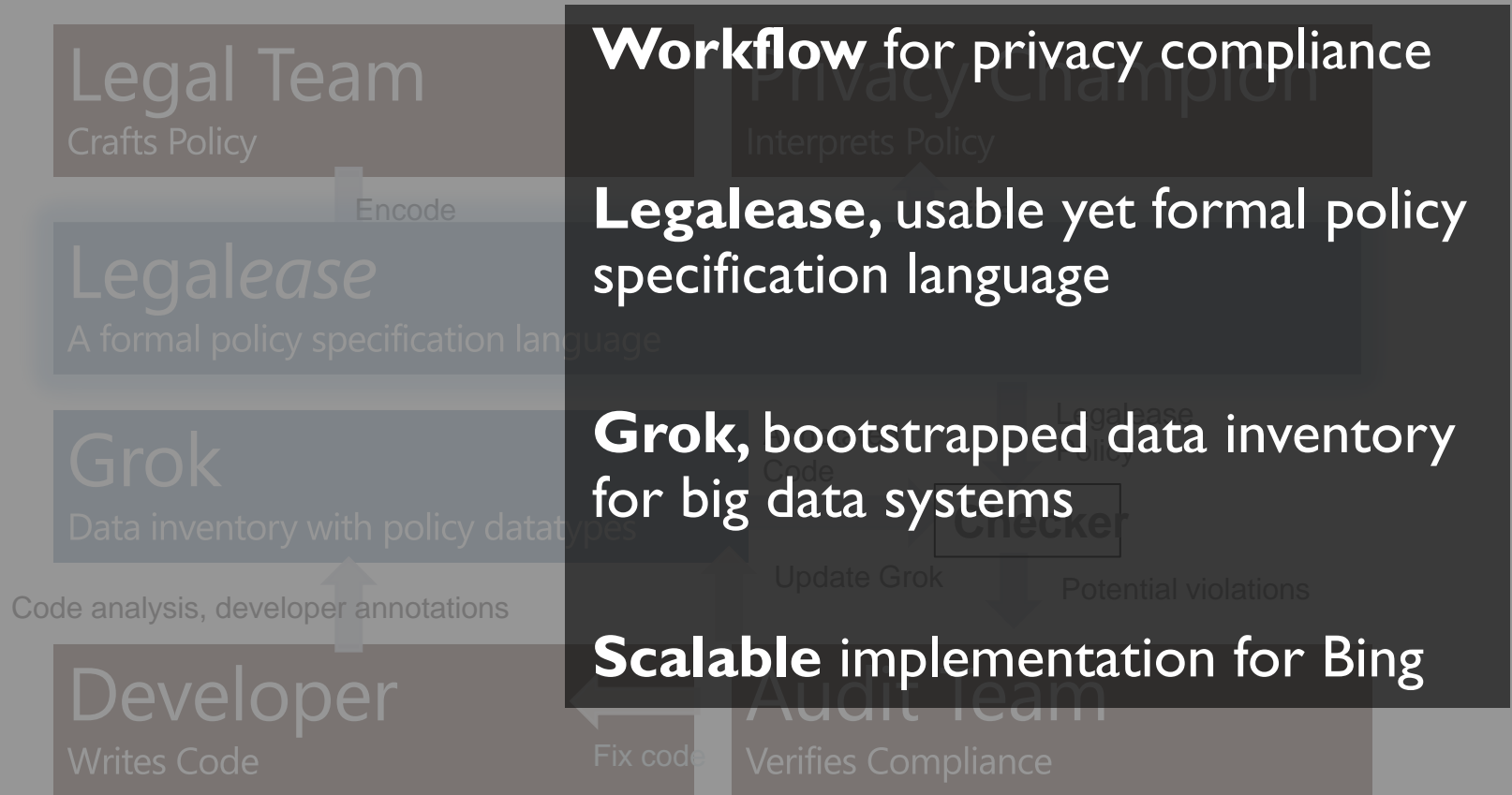


Fig. 9. Number of GROK data flow graph nodes added each day

- 77,000 jobs run each day
  - By 7000 entities
  - 300 functional groups
- 1.1 million unique lines of code
  - 21% changes on avg, daily
  - 46 million table schemas
  - 32 million files
- Manual audit infeasible
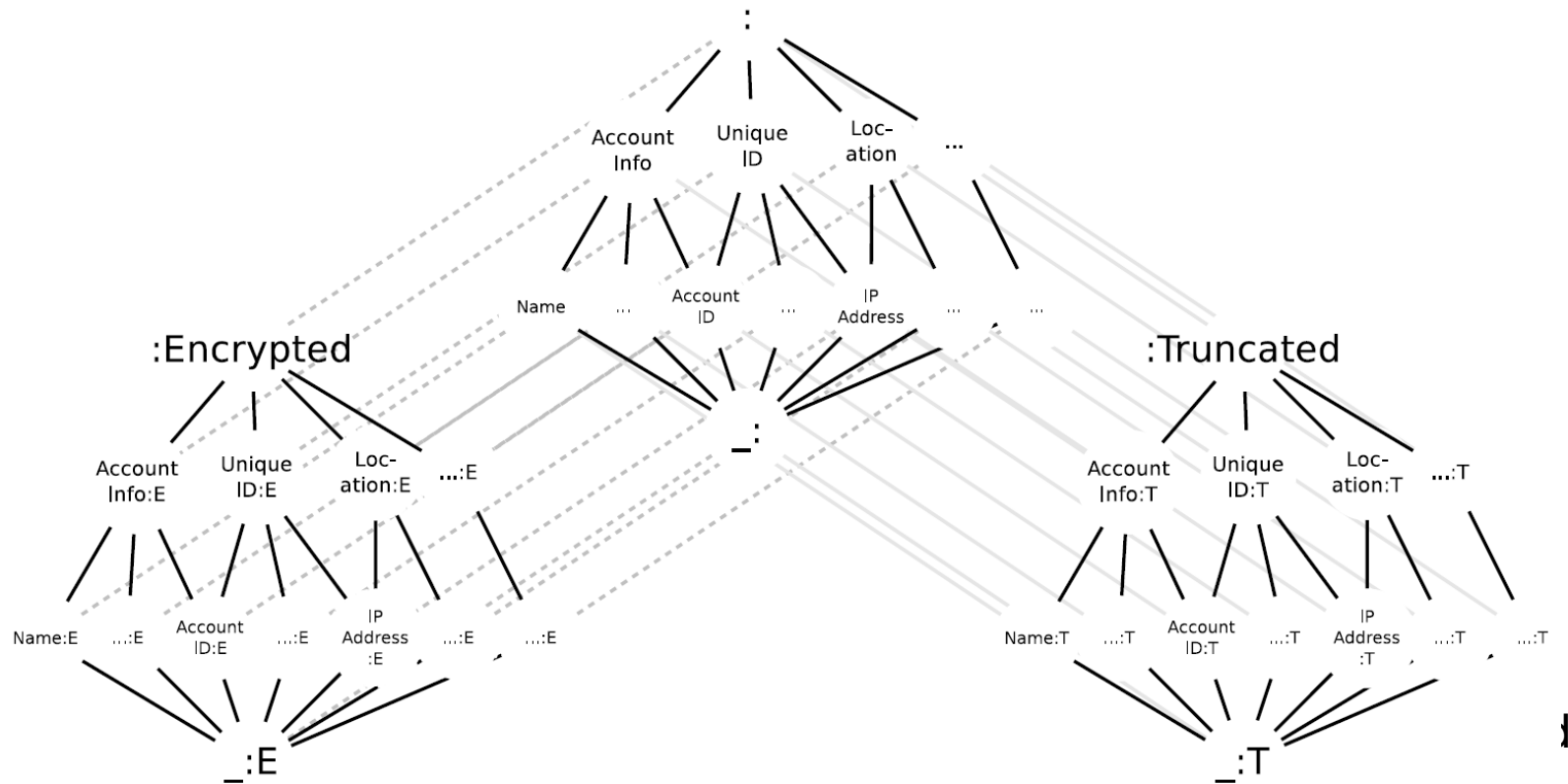- Information flow analysis takes ~30 mins daily

# A Streamlined Audit Workflow

# A Streamlined Audit Workflow

**Legal Team**
Crafts Policy

**Privacy Champion**
Interprets Policy

Encode

**Legal*ease***
A formal policy specification language

**Grok**
Data inventory with policy datatypes

Legalease
Policy

Code

Checker

Code analysis, developer annotations

Update Grok

Potential violations

**Developer**
Writes Code

Fix code

**Audit Team**
Verifies Compliance

---

**Workflow** for privacy compliance

**Legalease,** usable yet formal policy specification language

**Grok,** bootstrapped data inventory for big data systems
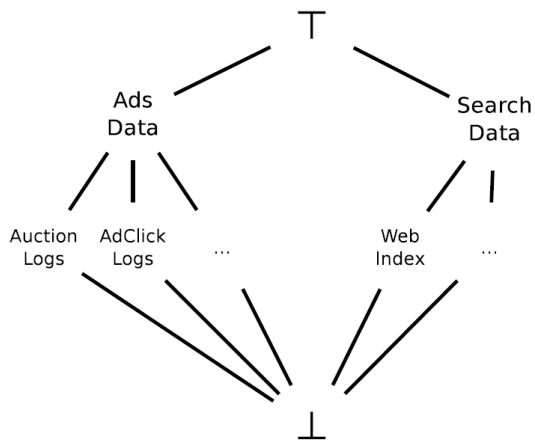
**Scalable** implementation for Bing

# Reference

- S. Sen, S. Guha, A. Datta, S. Rajamani, J. Tsai, J. M. Wing, Bootstrapping Privacy Compliance in Big Data Systems, in *Proceedings of 35th IEEE Symposium on Security and Privacy*, May 2014.
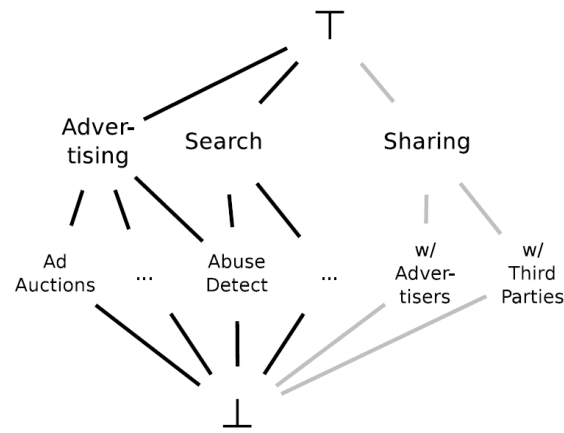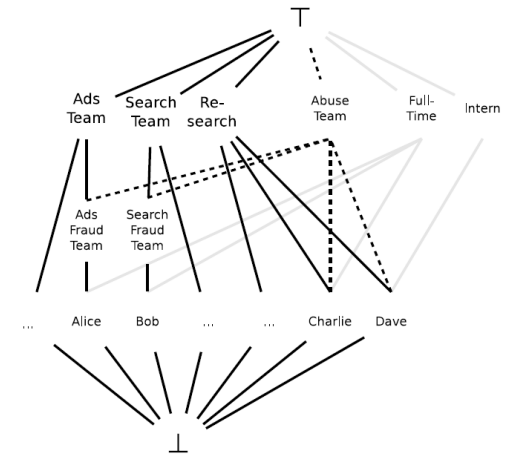
# Policy Labels : Datatypes

# Policy Types : Concept Lattices



*InStore* Lattice      *UseForPurpose* Lattice   *AccessByRole* Lattice

# Formal Semantics

$$\frac{\boxed{T^G \sqsubseteq T^C} \quad \exists_i D_i \ \text{denies} \ T^G}{\text{ALLOW} \ T^C \ \text{EXCEPT} \ D_1 \cdots D_m \ \text{denies} \ T^G} \ (\text{A}_2)$$

Based on Lattice Orderings on Policy Types

# Formal Semantics

$$\frac{T^G \sqsubseteq T^C \quad \boxed{\exists_i D_i \quad \text{denies} \quad T^G}}{\text{ALLOW } T^C \text{ EXCEPT } \boxed{D_1 \cdots D_m} \text{ denies } T^G} \quad (A_2)$$

**Recursively check exceptions**

ALLOW clauses have DENY clauses as exceptions

Top Level clause determines Blacklist/Whitelist

# Formal Semantics

$$\frac{T^G \sqsubseteq T^C \quad \exists_i D_i \ \text{denies} \ T^G}{\text{ALLOW} \ T^C \ \text{EXCEPT} \ D_1 \cdots D_m \ \text{denies} \ T^G} \ (\text{A}_2)$$

Structural properties about semantics

- Always returns an unambiguous answer
- Weakening lattice orderings makes policy more permissive

# Encoding

```
ALLOW
EXCEPT
    DENY DataType IPaddress:Expired
    DENY DataType UniqueIdentifier:Expired
    DENY DataType SearchQuery, PII InStore Store
    DENY DataType UniqueIdentifier, PII InStore Store

    DENY DataType BBEPData UseForPurpose Advertising


    DENY DataType BBEPData, PII InStore Store



    DENY DataType BBEPData:Expired

    DENY DataType UserProfile, PII InStore Store


    DENY DataType PII UseForPurpose Advertising
    DENY DataType PII InStore AdStore

    DENY DataType SearchQuery UseForPurpose Sharing
    EXCEPT
        ALLOW DataType SearchQuery:Scrubbed
```

◁ "we remove the entirety of the IP address after 6 months"

◁ "[we remove] cookies and other cross session identifiers, after 18 months"

◁ "We store search terms (and the cookie IDs associated with search terms) separately from any account information that directly identifies the user, such as name, e-mail address, or phone numbers."

◁ "we do not use any of the information collected through the Bing Bar Experience Improvement Program to identify, contact or target advertising to you"

◁ "we take steps to store [information collected through the Bing Bar Experience Improvement Program] separately from any account information we may have that directly identifies you, such as name, e-mail address, or phone numbers"

◁ "we delete the information collected through the Bing Bar Experience Program at eighteen months."

◁ "we store page views, clicks and search terms used for ad targeting separately from contact information you may have provided or other data that directly identifies you (such as your name, e-mail address, etc.)."

◁ "our advertising systems do not contain or use any information that can personally and directly identify you (such as your name, email address and phone number)."

◁ "Before we [share some search query data], we remove all unique identifiers such as IP addresses and cookie IDs from the data."