TaintStream: Fine-Grained Taint Tracking for Big Data Platforms through Dynamic Code Translation

Chengxu Yang, Yuanchun Li, Mengwei Xu, Zhenpeng Chen

Yunxin Liu, Gang Huang, Xuanzhe Liu







...



Privacy requirements (e.g. GDPR, HIPAA, COPPA, data provider's requirements, etc.) User data can be used only if:

- The requester has the permission
- Data is within retention period
- Data is properly anonymized
- The user doesn't choose opt-out



Privacy requirements (e.g. GDPR, HIPAA, COPPA, data provider's requirements, etc.)

Current practice:

User data can be used only if:

- The requester has the permission
- Data is within retention period
- Data is properly anonymized
- The user doesn't choose opt-out

Lots of manual code/data review

Case-by-case solution for each privacy requirement

Only support coarse-grained data management



Privacy requirements (e.g. GDPR, HIPAA, COPPA, data provider's requirements, etc.)

Current practice:

User data can be used only if:

- The requester has the permission
- Data is within retention period
- Data is properly anonymized
- The user doesn't choose opt-out

Problems:

Lots of manual code/data review

Case-by-case solution for each privacy requirement

Only support coarse-grained data management

Time-consuming

Error-prone

Data fragmentation & redundancy

A more systematic solution is needed!

An Example



- Sensitive information may be propagated *across various datasets*
- Data is transformed by *complicated scripts*
- Different compliance policies need to apply to all the datasets

Potential Solution: Taint Tracking



The original data is tainted with a tag

- e.g. variable c is tagged as "sensitive"
- Taint tags are propagated during program execution
 - e.g. *a* is tainted after a = b + c

out0 out0 out1

out1

VM go

v0 == v0 ta

v1 ==

v1 ta v2 ==

v4 tc

- Privacy policies can be enforced by checking the tags
 - e.g. *a* cannot be leaked if it's tagged as "sensitive"

Typical way to achieve taint tracking:

System modification

- Add a taint tag field to each register/variable
- Propagate taint tags in each operator [TaintDroid OSDI'10, LIFT MICRO'06]

		Op Format	Op Semantics	Taint Propagation
taint tag		const-op $v_A C$	$v_A \leftarrow C$	$\tau(v_A) \leftarrow \emptyset$
		move-op $v_A v_B$	$v_A \leftarrow v_B$	$\tau(v_A) \leftarrow \tau(v_B)$
Ladard Law		move-op-R vA	$v_A \leftarrow R$	$\tau(v_A) \leftarrow \tau(R)$
taint tag		return-op v_A	$R \leftarrow v_A$	$\tau(R) \leftarrow \tau(v_A)$
unused)		move-op-E vA	$v_A \leftarrow E$	$\tau(v_A) \leftarrow \tau(E)$
оор		throw-op vA	$E \leftarrow v_A$	$\tau(E) \leftarrow \tau(v_A)$
		unary-op vA vB	$v_A \leftarrow \otimes v_B$	$\tau(v_A) \leftarrow \tau(v_B)$
		binary-op $v_A v_B v_C$	$v_A \leftarrow v_B \otimes v_C$	$\tau(v_A) \leftarrow \tau(v_B) \cup \tau(v_C)$
= local0	•	binary-op $v_A v_B$	$v_A \leftarrow v_A \otimes v_B$	$\tau(v_A) \leftarrow \tau(v_A) \cup \tau(v_B)$
aint tag	•	binary-op $v_A v_B C$	$v_A \leftarrow v_B \otimes C$	$\tau(v_A) \leftarrow \tau(v_B)$
= local1		aput-op $v_A v_B v_C$	$v_B[v_C] \leftarrow v_A$	$\tau(v_B[\cdot]) \leftarrow \tau(v_B[\cdot]) \cup \tau(v_A)$
aint taa		aget-op $v_A v_B v_C$	$v_A \leftarrow v_B[v_C]$	$\tau(v_A) \leftarrow \tau(v_B[\cdot]) \cup \tau(v_C)$
attice cog		sput-op $v_A f_B$	$f_B \leftarrow v_A$	$\tau(f_B) \leftarrow \tau(v_A)$
= in0		sget-op $v_A f_B$	$v_A \leftarrow f_B$	$\tau(v_A) \leftarrow \tau(f_B)$
		iput-op $v_A v_B f_C$	$v_B(f_C) \leftarrow v_A$	$\tau(v_B(f_C)) \leftarrow \tau(v_A)$
aint tag		iget-op $v_A v_B f_C$	$v_A \leftarrow v_B(f_C)$	$\tau(v_A) \leftarrow \tau(v_B(f_C)) \cup \tau(v_B)$

Challenges

- Privacy compliance is not a set of precise computational rules
 - The taint tags should be flexible enough to support diverse privacy polices



- Big data platforms are complex and difficult to do system hacking
 - Runtime modification is unwanted for security and maintainability reasons
- Static approaches usually have oversimplified assumptions for real programs
 - Data schema is difficult to infer statically

System modification unwanted



Our approach (TaintStream)

Idea: *Statically rewrite* the script to let it *self-translate at runtime* and propagate *taint tags embedded in data frames*. Meanwhile, support diverse compliance requirements with *formal and flexible policy definition*.

Our approach (TaintStream)

Idea: *Statically rewrite* the script to let it *self-translate at runtime* and propagate *taint tags embedded in data frames*. Meanwhile, support diverse compliance requirements with *formal and flexible policy definition*.

- Create a parallel element to store taint tags for each data element
 The tag type can be flexibly customized according to the privacy policy
- Annotate the scripts through static code instrumentation
 - Program logic is changed by modifying the scripts instead of the system
- Complete translation at runtime and propagate taint tags on the fly
 Data schema is obtained in the dynamic context

Inspirations from source-to-source compilation approaches [JFlow POPL'99, TaintART CCS'16], while tailored for big data processing.

Workflow of TaintStream

Original script

conversationId	msg	date		
ххх	"Hi, I am"	10/3		
ууу	"Thanks"	12/4		

Original data

Flexible support of diverse privacy management tasks

Data Retention (GDPR)

The raw data and the data inferred from it must be deleted after a certain retention period (e.g., three months)

Access Control (organizational)

The access to certain data should be restricted if the requester is unauthorized.

User Data Erasure (GDPR)

When a user requests to be forgotten, the data collected from him/her and generated based on it must be deleted. **Tag type**: Integer *T*. # timestamp of the expiration date **Tag init**: : Set *T* to the time when the data expires. **Tag merge**: $Minimum(T_1, T_2)$ **Enforcement**: Periodically scan the datasets and delete any record *x* if the current timestamp is behind its timestamp T_x .

Tag type: Set of identifiers *S*. # set of authorized analysts **Tag init**: : $S_x \leftarrow$ the analysts with access to the raw data *x*. **Tag merge**: *Intersect*(S_1, S_2) **Enforcement**: A data record *x* is only visible to the analysts in its authorized analysts set S_x .

Tag type: Set of identifiers *S*. # set of involved data providers **Tag init**: : $S_x \leftarrow$ the provider the raw data *x* **Tag merge**: $Union(S_1, S_2)$ **Enforcement**: If a user *u* asks to be forgotten, delete any record *x* whose tag S_x contains *u*.

Dynamic translation

Index	Original	Translated					
1	$\phi(source)$	source _{tagged}					
2	$\phi(df.transformation)$	$\phi(df).\phi(transformation)$					
3	$\phi(select(Col))$	$select(\phi(Col))$					
4	$\phi(drop(Col))$	$drop(\phi(Col))$					
5	$\phi(orderBy(Col))$	$orderBy(V_{\phi(Col)})$					
6	$\phi(filter(Col))$	$filter(V_{\phi(Col)})$					
7	$\phi(join(df_2, on = Col))$	$join(\phi(df_2), on = V_{\phi(Col)})$					
8	$\phi(union(df_2))$	$union(\phi(df_2))$					
9	$\phi(withColumn(str, col))$	$withColumn(str, \phi(col))$					
	$\phi(groupBy(Col_1)$	$groupBy(V_{\phi(Col_1)})$					
10	$.agg(F_{agg_1}(Col_2),))$	$.agg(tagAgg(T_{\phi(Col_1)}), \phi(F_{agg_1}(Col_2),)$					
		$.map(x \to (\langle x_1, x_2 \rangle, x_3, x_4,))$					
	$\phi(map(x \rightarrow$	$map(x \rightarrow$					
(11)	$(F_i(x_{i_1}, x_{i_2},),$	$(\langle F_i(V_{x_{i_1}}, V_{x_{i_2}},), tagMerge(T_{x_{i_1}}, T_{x_{i_2}},)\rangle,$					
9	$F_j(x_{j_1}, x_{j_2},),$	$\langle F_j(V_{x_{j_1}}, V_{x_{j_2}},), tagMerge(T_{x_{j_1}}, T_{x_{j_2}},) \rangle,$					
)))))					
	$\phi(reduce(a, b \rightarrow$	$reduce(a, b \rightarrow$					
12	$(G_i(a_i,b_i),\ldots)))$	$(\langle G_i(V_{a_i}, V_{b_i}), tagMerge(T_{a_i}, T_{b_i})\rangle, \langle \rangle,)$					
13	$\phi(column_name))$	column_name					
14	$\phi(const)$	$\langle const, \perp \rangle$					
15	$\phi(Func(Col_1, Col_2,))$	$ \langle Func(V_{\phi(Col_1)}, V_{\phi(Col_2)},), \\ tagMerge(T_{\phi(Col_1)}, T_{\phi(Col_2)},) \rangle $					
\bigcirc							
Translation rules subset							

Parse the code at runtime and recursively apply the rules until no rule can apply.

Two design principles of translation rules:

- 1. (Non-interference) The translated code must not change the effect of the original code.
- 2. (**Conservativeness**) Should not tolerant false negatives.

Other designs:

- **Correctness guarantee**: Formally prove the noninterference property of the rules
- Exception handling: Snapshot and fallback to the conservative taint status on exceptions
- **Performance optimization**: Fuse non-conflicting taint propagation operations to save time

Demo: The datasets in TaintStream with taint tags

• Example: each taint tag is a Boolean value

UserId	Filename	DocumentTitle	Author		UserId		Filename	Docume	ntTitle		Author
System Account	a5eb6988-c9ad-44b	a5eb6988-c9ad-44b	System Account	[System Accou	rt,	[a5eb6988-c9	ad-44	[a5eb6988-c9a	d-44	System Acco	nt,
System Account	a5eb6988-c9ad-44b	a5eb6988-c9ad-44b	System Account	[System Accou	nt,	[a5eb6988-c9	ad-44	[a5eb6988-c9a	d-44	System Acco	nt,
温 盈盈	20200907-yingying	The Upper Bound o	温 盈盈	[温 盈盈]	true]	[20200907-yi	ngyin	The Upper Bo	und]	[温 盈盈	, true]
	学校同意函.pdf	学校同意函		l Ī	true]	「学校同意函.pdf	, false]	[学校同意函,	false]		, true]
	学校同意函.pdf	学校同意函		Ĩ Ĩ	true]]	「学校同意函.pdf	, false]	「学校同意函,	false]]		, true]
	学校同意函.pdf	学校同意函		Ì Ì	true][「学校同意函.pdf	, false]	「学校同意函,	false]		, true]
	健康宝.png	健康宝		Î Î	true]]	[健康宝.png	, false]	[健康宝,	false]]		, true]
	健康宝.png	健康宝		Ì Ì	true][[健康宝.png	, false]	「健康宝,	false][, true]
	健康宝.png	健康宝	i	Ì Ì	true	[健康宝.png	, false]	[健康宝,	false]]		, true]
i	自愿返岗承诺书.pdf	自愿返岗承诺书	i	i ī	trueii	自愿返岗承诺书.pdf	, false]	[自愿返岗承诺书,	false]		, true]
	自愿返岗承诺书.pdf	自愿返岗承诺书		i i	true	自愿返岗承诺书.pdf	, false]	[自愿返岗承诺书,	false		, true]
	自愿返岗承诺书.pdf	自愿返岗承诺书	i	i ī	true	自愿返岗承诺书.pdf	, false]	自愿返岗承诺书,	false]		, true]
Chengxu Yang (FA	学校同意函.pdf	学校同意函	Chengxu Yang (FA)	[Chengxu Yang	(FA	[学校同意函.pdf	, false]	[学校同意函]	false]]	Chengxu Yan	(FA
Chengxu Yang (FA	学校同意函.pdf	学校同意函	Chengxu Yang (FA)	[Chengxu Yang	(FA	「学校同意函.pdf	, false]	「学校同意函,	false]	Chengxu Yan	(FA
Chengxu Yang (FA	学校同意函.pdf	学校同意函	Chengxu Yang (FA)	[Chengxu Yang	(FA	「学校同意函.pdf	, false]	「学校同意函,	false]]	Chengxu Yan	(FA
Chengxu Yang (FA	学校同意函.pdf	学校同意函	Chengxu Yang (FA	[Chengxu Yang	(FA	「学校同意函.pdf	, false]	「学校同意函,	false]	Chengxu Yan	(FA
Chengxu Yang (FA	学校同意函.pdf	学校同意函	Chengxu Yang (FA)	[Chengxu Yang	(FA	「学校同意函.pdf	, false]	「学校同意函,	false]]	Chengxu Yan	(FA
Chengxu Yang (FA	学校同意函.pdf	学校同意函	Chengxu Yang (FA	[Chengxu Yang	(FA)	「学校同意函.pdf	, false]	「学校同意函,	false]	Chengxu Yan	(FA)
Chengxu Yang (FA	学校同意函.pdf	学校同意函	Chengxu Yang (FA	[Chengxu Yang	(FA)	「学校同意函.pdf	, false]	「学校同意函,	false]]	Chengxu Yan	(FA
Chengxu Yang (FA	自愿返岗承诺书.pdf	自愿返岗承诺书	Chengxu Yang (FA	[Chengxu Yang	(FA	自愿返岗承诺书.pdf	, false]	[自愿返岗承诺书,	false]]	Chengxu Yan	(FA
+		·	+	+	+-	<u></u>					
only showing top 20 rd	DWS			only showing to	p 20 row	N S					

The output dataset of the original script

The output dataset when running with TaintStream

Experiment setup

- Query scripts: 7 real-world scripts and 33 self-built scripts
- Baseline: PlanAnlyzer, PySa (SOTA static analyzers)
- Environment: 4-node cluster
- Metrics:
 - precision, recall
 - running time, storage overhead

Accuracy

93.0% precision and 100% recall

	Column Level			Cell Level				
Script Name	PySa PlanAnalyzer		TaintStream	TaintStream				
		correct/tota	1	precision	recall			
Basic								
orderBy_1	4/6	6/6	6/6	100.0%	100.0%			
orderBy_2	1/3	3/3	3/3	100.0%	100.0%			
select_1	3/6	6/6	6/6	100.0%	100.0%			
select_2	1/2	2/2	2/2	100.0%	100.0%			
withColumn_1	3/3	3/3	3/3	100.0%	100.0%			
withColumn_2	3/5	5/5	5/5	100.0%	100.0%			
withColumn_3	3/3	3/3	3/3	100.0%	100.0%			
		GroupBy	,					
count_1	1/3	3/3	3/3	100.0%	100.0%			
count_2	2/3	3/3	3/3	100.0%	100.0%			
count_3	2/5	5/5	5/5	100.0%	100.0%			
count_4	2/4	4/4	4/4	100.0%	100.0%			
statistics_1	4/5	5/5	5/5	100.0%	100.0%			
statistics_2	3/5	5/5	5/5	100.0%	100.0%			
Join								
inner_join_1	2/4	4/4	4/4	100.0%	100.0%			
inner_join_2	3/5	5/5	5/5	100.0%	100.0%			
inner_join_3	2/5	5/5	5/5	100.0%	100.0%			
left_join	3/5	5/5	5/5	100.0%	100.0%			
right_join	3/5	5/5	5/5	100.0%	100.0%			
outer join	3/5	5/5	5/5	100.0%	100.0%			

Structured Data							
array_type_1 *	4/5	5/5	4/5	80.0%	100.0%		
array_type_2	2/3	3/3	3/3	100.0%	100.0%		
struct_type_1	2/5	1/5	4/5	100.0%	100.0%		
struct_type_2	2/3	2/3	3/3	100.0%	100.0%		
<pre>struct_type_3 *</pre>	3/4	2/4	3/4	66.7%	100.0%		
		UDF					
udf_1	1/4	4/4	4/4	100.0%	100.0%		
udf_2	1/1	1/1	1/1	100.0%	100.0%		
udf_3	0/1	0/1	0/1	0.0%	100.0%		
class_udf_1	2/2	2/2	2/2	100.0%	100.0%		
class_udf_2	1/1	1/1	1/1	100.0%	100.0%		
		Map Reduce					
map_reduce_1	0/2	0/2	2/2	100.0%	100.0%		
map_reduce_2	1/4	0/4	4/4	100.0%	100.0%		
map_reduce_3	1/4	0/4	4/4	100.0%	100.0%		
map_reduce_4 *	1/4	0/4	1/4	21.9%	100.0%		
Summary	69/125	103/125	118/125	93.0%	100.0%		

System Overhead



Script Name	Original Output Size (Byte)	Tag Size	Overhead
email tokenizer	2,697,017	19,952 <i>x</i>	0.74%
never replied emails	246,129	5,000 <i>x</i>	2.03%
top name	16,229	153 <i>x</i>	0.94%
action provider	32,073	636x	1.98%
email provider	10,246,702	44,275 <i>x</i>	0.43%
data statistics	633	25x	3.95%
extract documents	1,097	48x	4.38%
top name action provider email provider data statistics extract documents	16,229 32,073 10,246,702 633 1,097	153x 636x 44,275x 25x 48x	0.94% 1.98% 0.43% 3.95% 4.38%

🗑 12.7% running time overhead



Take Away

- Light-weight fine-grained taint tracking for big data platforms
- Various privacy management tasks with little overhead
- With no changes to the system and developers

Data and framework: https://github.com/PrivacyStreams/TaintStream



Thanks!





