

TaintStream Appendix

A DYNAMIC TRANSLATION RULES

Table 1: Examples of dynamic code translation rules. $\phi(\text{code})$ represents applying translation rules to the code. $\langle \text{value}, \text{tag} \rangle$ is to pack a value field and a taint tag field into a structured value field (the new field name is set to the same as the value field name). V_{Col} and T_{Col} are used to get the value and tag of the column respectively. $\text{tagMerge}(\text{tag}_1, \text{tag}_2, \dots)$ and $\text{tagAgg}(\text{tag})$ are to merge multiple taint tags or aggregate a group of tags using a customized merging function. \perp represents the default tag for const/insensitive values.

Index	Original	Translated
①	$\phi(\text{source})$	$\text{source}_{\text{tagged}}$
②	$\phi(df.\text{transformation})$	$\phi(df).\phi(\text{transformation})$
③	$\phi(\text{select}(Col))$	$\text{select}(\phi(Col))$
④	$\phi(\text{drop}(Col))$	$\text{drop}(\phi(Col))$
⑤	$\phi(\text{orderBy}(Col))$	$\text{orderBy}(V_{\phi(Col)})$
⑥	$\phi(\text{filter}(Col))$	$\text{filter}(V_{\phi(Col)})$
⑦	$\phi(\text{join}(df_2, \text{on} = Col))$	$\text{join}(\phi(df_2), \text{on} = V_{\phi(Col)})$
⑧	$\phi(\text{union}(df_2))$	$\text{union}(\phi(df_2))$
⑨	$\phi(\text{withColumn}(\text{str}, col))$	$\text{withColumn}(\text{str}, \phi(col))$
⑩	$\phi(\text{groupBy}(Col_1))$	$\text{groupBy}(V_{\phi(Col_1)})$
⑪	$\phi(\text{agg}(\text{Func}(Col_2, \dots)))$	$\text{agg}(\text{tagAgg}(T_{\phi(Col_1)}), \phi(\text{Func}(Col_2, \dots)))$
⑫	$\phi(\text{map}(x \rightarrow (F_i(x_{i_1}, x_{i_2}, \dots), F_j(x_{j_1}, x_{j_2}, \dots), \dots)))$	$\text{map}(x \rightarrow ((F_i(V_{x_{i_1}}, V_{x_{i_2}}, \dots), \text{tagMerge}(T_{x_{i_1}}, T_{x_{i_2}}, \dots)), (F_j(V_{x_{j_1}}, V_{x_{j_2}}, \dots), \text{tagMerge}(T_{x_{j_1}}, T_{x_{j_2}}, \dots), \dots)))$
⑬	$\phi(\text{reduce}(a, b \rightarrow (G_i(a_i, b_i), \dots)))$	$\text{reduce}(a, b \rightarrow ((G_i(V_{a_i}, V_{b_i}), \text{tagMerge}(T_{a_i}, T_{b_i})), (\dots), \dots))$
⑭	$\phi(\text{column_name})$	column_name
⑮	$\phi(\text{const})$	$\langle \text{const}, \perp \rangle$
⑯	$\phi(\text{Func}(Col_1, Col_2, \dots))$	$\langle \text{Func}(V_{\phi(Col_1)}, V_{\phi(Col_2)}, \dots), \text{tagMerge}(T_{\phi(Col_1)}, T_{\phi(Col_2)}, \dots) \rangle$
⑰	$\phi(\text{Fagg}(Col))$	$\langle \text{Fagg}(V_{\phi(Col)}), \text{tagAgg}(T_{\phi(Col)}) \rangle$

B PROOF OF CORRECTNESS

We take the *withColumn* operation as an example to prove the correctness of our translation function ϕ , i.e., the output dataframe produced by the translated *withColumn* is equivalent to the dataframe produced by the original *withColumn*, excluding the taint tags.

We first formulate some basic properties of dataframes and columns, which are regarded as the axioms in our proof.

Dataset properties. Suppose df , c , r represent a dataframe, a column, and a row respectively, we have:

- (1) $df[c][r]$ represents getting the value at column c and row r .
- (2) $df_1 = df_2 \Leftrightarrow \forall c, r : df_1[c][r] = df_2[c][r]$

Taint operations. Suppose $df^t = \text{taint}(df)$, we have:

- (1) $df^t[c][r] = \langle df[c][r], \text{tag} \rangle$.
- (2) $df^t[c][r].\text{val} = df[c][r]$.

Definition of *withColumn*. Suppose $df_2 = df_1.\text{withColumn}(\text{str}, Col)$, for \forall column c and row $r \in df_2$, we have:

- (1) $df_2[c][r] = df_1[c][r]$, if c 's column name is not str ;

- (2) $df_2[c] = Col$, if c 's column name is str .

We then prove the following lemma to show the correctness of the translation rules on the column, i.e., Rule ⑬ ~ ⑰.

LEMMA B.1. For $\forall df$, $df^t = \text{taint}(df)$, then for \forall column c and row r , $df^t[\phi(c)][r].\text{val} = df[c][r]$.

We prove the lemma using mathematical induction:

- (1) if c is column_name, getting a column by its name will return a tainted column, so $df^t[\phi(c)][r].\text{val} = df[c][r]$.
- (2) else if c is const, according to the Rule ⑭, $\phi(c) = \langle \text{const}, \perp \rangle$, so $df^t[\phi(c)][r].\text{val} = df[c][r]$.
- (3) else if c is $\text{Func}(c_1, c_2, \dots)$, we have:

$$\begin{aligned}
 \phi(c) &= \phi(\text{Func}(c_1, c_2, \dots)) \\
 &= \langle \text{Func}(V_{\phi(c_1)}, V_{\phi(c_2)}, \dots), \text{tagMerge}(\text{tag}_1, \text{tag}_2, \dots) \rangle \\
 &\dots\dots\dots \text{Rule ⑮} \\
 &= \langle \text{Func}(V_{\phi(c_1)}, V_{\phi(c_2)}, \dots), \text{tag} \rangle \dots\dots\dots \text{tag merge} \\
 &= \langle \text{Func}(c_1, c_2, \dots), \text{tag} \rangle \dots\dots\dots \text{mathematical induction} \\
 &= \langle c, \text{tag} \rangle \dots\dots\dots c \text{ is } \text{Func}(c_1, c_2, \dots)
 \end{aligned}$$

so $df^t[\phi(c)][r].\text{val} = df^t[\langle c, \text{tag} \rangle][r].\text{val} = df[c][r]$

- (4) else if c is $\text{Fagg}(c)$, the proof is similar to (3).

Overall, we have proved $df^t[\phi(c)][r].\text{val} = df[c][r]$.

Now we prove the correctness of the translation rule of *withColumn*. Specifically, we need to prove that the output dataframe through the translated *withColumn* operation is equivalent to the dataframe through the original *withColumn* operation, excluding the taint tags, which can be formalized as follows.

PROOF. For $\forall df_1$, $df_1^t = \text{taint}(df_1)$; and for $\forall \text{str}, Col$, $df_2 = df_1.\text{withColumn}(\text{str}, Col)$. Then for $\forall c, r$, $\phi(df_2)[c][r].\text{val} = df_2[c][r]$.

$$\begin{aligned}
 \phi(df_2) &= \phi(df_1.\text{withColumn}(\text{str}, Col)) \\
 &= \phi(df_1).\phi(\text{withColumn}(\text{str}, Col)) \dots\dots\dots \text{Rule ②} \\
 &= df_1^t.\text{withColumn}(\text{str}, \phi(Col)) \dots\dots \text{Rule ①, Rule ⑨} \\
 &= df_1^t.\text{withColumn}(\text{str}, \langle Col, \text{tag} \rangle) \dots\dots\dots \text{Lamma B.1}
 \end{aligned}$$

If c 's column name is str , then

$$\phi(df_2)[c][r].\text{val} = \langle Col, \text{tag} \rangle.\text{val} = Col[r] = df_2[c][r]$$

Else, c 's column name is not str , then

$$\phi(df_2)[c][r].\text{val} = df_1^t[c][r].\text{val} = df_1[c][r] = df_2[c][r]$$

Overall, we have prove $\phi(df_2)[c][r].\text{val} = df_2[c][r]$ \square