# A Formal Model of Provenance in Distributed Systems

Issam Souilah[2]    Adrian Francalanza[1]    Vladimiro Sassone[2]

[1]Department of Computer Science
ICT, University of Malta

[2]DSSE, ECS
Southampton University

San Francisco, USA. February 2009

# Outline

Trust In a Distributed System

# Motivation

Trust In a Distributed System

- Distribution $\Rightarrow$ inherent **parallelism**.

Trust In a Distributed System

- Distribution $\Rightarrow$ inherent **parallelism**.
- Distribution $\Rightarrow$ no shared memory *i.e.,* **message passing**.

# Motivation

## Trust In a Distributed System

- Distribution $\Rightarrow$ inherent **parallelism**.
- Distribution $\Rightarrow$ no shared memory *i.e.,* **message passing**.
- Distribution $\Rightarrow$ lack of centralised coordination *i.e.,* **non-determinism**.

### Trust In a Distributed System

- Distribution $\Rightarrow$ inherent **parallelism**.
- Distribution $\Rightarrow$ no shared memory *i.e.,* **message passing**.
- Distribution $\Rightarrow$ lack of centralised coordination *i.e.,*
  **non-determinism**.

### Appropriate Calculus?
The **piCalculus** :

# Motivation

## Trust In a Distributed System

- Distribution $\Rightarrow$ inherent **parallelism**.
- Distribution $\Rightarrow$ no shared memory *i.e.,* **message passing**.
- Distribution $\Rightarrow$ lack of centralised coordination *i.e.,* **non-determinism**.

## Appropriate Calculus?

The **piCalculus** :

- Captures the characteristic features of our domain of study.

# Motivation

## Trust In a Distributed System

- Distribution $\Rightarrow$ inherent **parallelism**.
- Distribution $\Rightarrow$ no shared memory *i.e.,* **message passing**.
- Distribution $\Rightarrow$ lack of centralised coordination *i.e.,* **non-determinism**.

## Appropriate Calculus?

The **piCalculus** :

- Captures the characteristic features of our domain of study.
- Well-studied.

# Motivation

## Trust In a Distributed System

- Distribution $\Rightarrow$ inherent **parallelism**.
- Distribution $\Rightarrow$ no shared memory *i.e.,* **message passing**.
- Distribution $\Rightarrow$ lack of centralised coordination *i.e.,* **non-determinism**.

## Appropriate Calculus?

The **piCalculus** :

- Captures the characteristic features of our domain of study.
- Well-studied.
- Close connections to linear logic and resources
- . . .

# piCalculus Primer

## Names

$$a, b, \ldots \in \textsc{Names}$$

# piCalculus Primer

Names

$$a, b, \ldots \in \text{NAMES}$$

- denote points of interaction (rendez-vous **channels**)...

# piCalculus Primer

Names

$$a, b, \ldots \in \text{Names}$$

- denote points of interaction (rendez-vous **channels**)...
- ...and **values** which are transmitted during communication.

# piCalculus Primer

Processes

$$P \quad \| \quad Q \quad \| \quad R$$

# piCalculus Primer

Processes

$$a!b \quad \| \quad Q \quad \| \quad R$$

# piCalculus Primer

Processes

$$a!b \quad \| \quad a?x.c!x \quad \| \quad R$$

# piCalculus Primer

Processes

$$a!b \quad \| \quad a?x.c!x \quad \| \quad R$$

# piCalculus Primer

Processes

$$a!b \quad \| \quad a?x.x!c \quad \| \quad R$$

# piCalculus Primer

Processes

$$a!b \quad \| \quad a?x.x!c \quad \| \quad R$$

react on $a$

Processes

$$b!c \quad \| \quad R$$

# piCalculus Primer

Processes

$$b!c \quad \| \quad b?y.R'$$

# piCalculus Extension: Distribution!

From Processes to Systems

$$P \quad \| \quad Q \quad \| \quad R$$

# piCalculus Extension: Distribution!

From Processes to Systems

$$p, q, r, \ldots \in \text{TRUSTPRINCIPALS}$$

$$p[P] \quad \| \quad q[Q] \quad \| \quad p[R]$$

# piCalculus Extension: Distribution!

From Processes to Systems

$$p[a?x.P] \quad \| \quad q[a!v_1] \quad \| \quad r[R]$$

# piCalculus Extension: Distribution!

From Processes to Systems

$$p[a?x.P] \quad \| \quad q[a!v_1] \quad \| \quad r[R]$$

**across** units of trust

# piCalculus Extension: Distribution!

From Processes to Systems

$$p[a?x.P] \quad \| \quad q[a!v_1] \quad \| \quad r[a!v_2]$$

$\underbrace{\hspace{8cm}}$

**Market** of values!

# piCalculus Extension: Distribution!

From Processes to Systems

$$\overbrace{p[a?x.P] \quad \| \quad q[a!v_1]}^{\text{communication}} \quad \| \quad r[a!v_2]$$

$$\underbrace{p[a?x.P] \quad \| \quad q[a!v_1] \quad \| \quad r[a!v_2]}_{\textbf{Market of values!}}$$

# piCalculus Extension: Distribution!

From Processes to Systems

$$\overbrace{p[a?x.P] \quad \| \quad q[a!v_1] \quad \| \quad r[a!v_2]}^{\text{communication}}$$

**Market** of values!

What are the control mechanism need to assist consumers of data?

*How do we automate decisions based on trust?*

What are the control mechanism need to assist consumers of data?

*How do we automate decisions based on trust?*

- ▶ Static Analysis *(not scalable)*

What are the control mechanism need to assist consumers of data?

*How do we automate decisions based on trust?*

- Static Analysis *(not scalable)*
- Dynamic Analysis:
    - decisions need to be computationally lightweight.

    - decision criteria produced in lightweight fashion.

What are the control mechanism need to assist consumers of data?

*How do we automate decisions based on trust?*

- Static Analysis *(not scalable)*
- Dynamic Analysis:
  - decisions need to be computationally lightweight. *(Full-blown **verification methods** do not cut it!)*
  - decision criteria produced in lightweight fashion.

What are the control mechanism need to assist consumers of data?

*How do we automate decisions based on trust?*

- Static Analysis *(not scalable)*
- Dynamic Analysis:
    - decisions need to be computationally lightweight. *(Full-blown **verification methods** do not cut it!)*
    - decision criteria produced in lightweight fashion. *(**Proof-Carrying Code** does not cut it!)*

# Outline

Annotated Values

$$v :$$

Annotated Values

$$v : \kappa$$

# Provenance

### Annotated Values

$$v : \kappa$$

$$
\begin{aligned}
\kappa ::= &\ \epsilon & \text{empty provenance} \\
| &\ \alpha; \kappa & \text{sequenced provenace} \\
\alpha ::= &\ \mathbf{rcv}(p, \kappa) & \text{recieve action} \\
| &\ \mathbf{snd}(p, \kappa) & \text{send action}
\end{aligned}
$$

Annotated Values

$$p[a!v]$$

Annotated Values

$$p[a : \kappa_a ! v : \kappa_v]$$

# Provenance

## Annotated Values

$$p[a : \kappa_a \,!\, v : \kappa_v] \quad \| \quad q[a : \kappa_a' \,!\, v' : \kappa_{v'}.] \quad \| \quad p[a : \kappa_a'' \,!\, v'' : \kappa_{v''}.]$$

provenance is linear!

# Provenance Tracking

Automated:

Automated:

1. orthogonal to programming *(can be abstracted away)*

# Provenance Tracking

Automated:

1. orthogonal to programming *(can be abstracted away)*
2. ensures provenance annotation standardization.

# Provenance Tracking

Automated:

1. orthogonal to programming *(can be abstracted away)*
2. ensures provenance annotation standardization.
3. avoids circular reasoning with respect to trust.

# Provenance Tracking

Automated:

1. orthogonal to programming *(can be abstracted away)*
2. ensures provenance annotation standardization.
3. avoids circular reasoning with respect to trust.

Two tiered architecture:

Automated:

1. orthogonal to programming *(can be abstracted away)*
2. ensures provenance annotation standardization.
3. avoids circular reasoning with respect to trust.

Two tiered architecture:

- ► **Computation Layer:** describes computation of values and processes.

# Provenance Tracking

Automated:

1. orthogonal to programming *(can be abstracted away)*
2. ensures provenance annotation standardization.
3. avoids circular reasoning with respect to trust.

Two tiered architecture:

- **Computation Layer:** describes computation of values and processes.
- **Provenance Tracking Layer:** describes the aggregation of provenance information attached to data *(typically assigned to a trusted middleware)*

Operational Semantics

$$p[a!v] \parallel Q \quad \longrightarrow \quad a\langle v \rangle \parallel Q$$

Operational Semantics

$$p[a!v] \parallel Q \quad \longrightarrow \quad a\langle v \rangle \parallel Q$$

loose immediate provenance information!

# Provenance Tracking Semantics

*(Provenance Tracking)* Operational Semantics

$$p[a : \kappa_a ! v : \kappa_v] \parallel Q \quad \longrightarrow \quad a\langle v : \mathbf{snd}(p, \kappa_a); \kappa_v \rangle \parallel Q$$

provenance aggregation

Not-Automated!

Not-Automated!

- program with it to control non-derminism...

# Provenance Usage

Not-Automated!

- program with it to control non-derminism...
- ...using intuitive programming idioms/constructs

# Provenance Usage

Not-Automated!

- program with it to control non-derminism. . .
- . . . using intuitive programming idioms/constructs

Operational Semantics

$$a\langle v \rangle \parallel q[a?(x).Q] \longrightarrow q[Q\{^v\!/\!x\}]$$

# Provenance Usage

Not-Automated!

- program with it to control non-derminism. . .
- . . . using intuitive programming idioms/constructs

## Operational Semantics

$$a\langle v : \kappa_v \rangle \parallel q[a : \kappa_a ?(x \; \textbf{from} \; \pi).Q] \; \longrightarrow \; q[Q\{^v/x\}] \quad \text{if } \kappa_v \models \pi$$

$$\underbrace{\phantom{q[a : \kappa_a ?(x \; \textbf{from} \; \pi).Q]}}_{\text{provenance pattern matching}}$$

# Provenance Usage

Not-Automated!

- ▶ program with it to control non-derminism. . .
- ▶ . . . using intuitive programming idioms/constructs

## Operational Semantics

$$a\langle v : \kappa_v \rangle \parallel q[a : \kappa_a ?(x \textbf{ from } \pi).Q] \; \longrightarrow \; q[Q\{^{v : \textbf{rcv}(q, \kappa_a) \, ; \, \kappa_v}/x\}] \quad \text{if } \kappa_v \models \pi$$

$$\underbrace{\phantom{v : \textbf{rcv}(q, \kappa_a) \, ; \, \kappa_v}}_{\text{provenance aggregation}}$$

# Provenance Usage Example

### Client/Server

    *srv*   Name of server

    *ret*   Name of return channel on which server returns answer

$$Client = p[\,srv!\langle ret\rangle\,] \quad \| \quad p[\,ret?(x\ \textbf{from}\ \pi).P]$$
$$Server = q[\,srv?(y\ \textbf{from}\ *).y!\langle v\rangle\,]$$

## Client/Server

    *srv*   Name of server

    *ret*   Name of return channel on which server returns answer

$$Client = p[\,srv!\langle ret\rangle\,] \quad \| \quad p[\,ret?(x \textbf{ from } \pi).P]$$
$$Server = q[\,srv?(y \textbf{ from } *).y!\langle v\rangle\,]$$

# Provenance Usage Example

## Client/Server

  *srv*   Name of server

  *ret*   Name of return channel on which server returns answer

$$Client = p[\,srv!\langle ret\rangle\,] \quad \| \quad p[\,ret?(x\,\textbf{from}\,\pi).P\,]$$
$$Server = q[\,srv?(y\,\textbf{from}\,*).y!\langle v\rangle\,]$$

Client/Server

> *srv*    Name of server
>
> *ret*    Name of return channel on which server returns answer

$$Client = p[\, srv!\langle ret\rangle\,] \quad \| \quad p[\, ret?(x \textbf{ from } \pi).P\,]$$
$$Server = q[\, srv?(y \textbf{ from } *).y!\langle v\rangle\,]$$

# Provenance Usage Example

Client/Server

> *srv*   Name of server
>
> *ret*   Name of return channel on which server returns answer

$$Client = p[\, srv : \kappa^1_{srv} \, !\langle\, ret : \epsilon\,\rangle\,] \quad \| \quad p[\, ret : \epsilon \, ?(x \textbf{ from } \pi).P\,]$$
$$Server = q[\, srv : \kappa^2_{srv} \, ?(y \textbf{ from } *).y!\langle v : \kappa_v \rangle\,]$$

# Provenance Usage Example

## Client/Server

     *srv*   Name of server

     *ret*   Name of return channel on which server returns answer

$$Client = p[\, srv\!:\!\kappa_{srv}^{1}\,!\langle\, ret\!:\!\epsilon\rangle] \quad \| \quad p[\, ret\!:\!\epsilon\,?(x\,\textbf{from}\,\pi).P]$$

$$Server = q[\, srv\!:\!\kappa_{srv}^{2}\,?(y\,\textbf{from}\,*).y!\langle v\!:\!\kappa_{v}\rangle]$$

Client/Server

$$srv \quad \text{Name of server}$$
$$ret \quad \text{Name of return channel on which server returns answer}$$

$$Client = p[\, srv : \kappa_{srv}^1 \,! \langle \, ret : \epsilon \rangle] \quad \| \quad p[\, ret : \epsilon \,?(x \textbf{ from } \pi).P]$$
$$Server = q[\, srv : \kappa_{srv}^2 \,?(y \textbf{ from } *).y! \langle v : \kappa_v \rangle]$$

$$\pi = \textbf{snd}(q, \textbf{rcv}(q, \textbf{snd}(p, \epsilon))); *$$

# Outline

Correctness Intuition

- provenance attached to values records **history related** to that **value**.
- provenance of a value is correct if it describes a **partial history** which corresponds to the **total history** of events.

# Provenance Correctness

History represented by Logs

$$\phi ::= \emptyset \quad | \quad \rho; \phi \qquad \qquad \text{logs}$$

$$\rho ::= \textbf{rcv}(p) \quad | \quad \textbf{snd}(p) \qquad \text{log actions}$$

# Provenance Correctness

## History represented by Logs

$$\phi ::= \emptyset \quad | \quad \rho; \phi \qquad \qquad \text{logs}$$
$$\rho ::= \textbf{rcv}(p) \quad | \quad \textbf{snd}(p) \qquad \text{log actions}$$

## Sub-Log Comparison

$$\text{cmp1} \frac{}{\emptyset \leq \phi} \qquad \text{cmp2} \frac{\phi \leq \phi'}{\rho; \phi \leq \rho; \phi'} \qquad \text{cmp3} \frac{\phi \leq \phi'}{\phi \leq \rho; \phi'}$$

Monitored Systems

$$\phi \triangleright p[a : \kappa_a ! v : \kappa_v] \parallel Q \quad \longrightarrow_{\mathsf{mon}} \quad \mathbf{snd}(p); \phi \triangleright a\langle v : \mathbf{snd}(p, \kappa_a); \kappa_v\rangle \parallel Q$$

## Monitored Systems

$$\phi \triangleright p[a : \kappa_a ! v : \kappa_v] \parallel Q \quad \longrightarrow_{\text{mon}} \quad \textbf{snd}(p); \phi \triangleright a\langle v : \textbf{snd}(p, \kappa_a); \kappa_v\rangle \parallel Q$$

**Erasure Function:**

$$|-| : \text{MONITOREDSYS} \to \text{SYS}$$

# Provenance Correctness

## Monitored Systems

$$\phi \triangleright p[a : \kappa_a ! v : \kappa_v] \parallel Q \quad \longrightarrow_{\mathsf{mon}} \quad \mathbf{snd}(p); \phi \triangleright a\langle v : \mathbf{snd}(p, \kappa_a); \kappa_v\rangle \parallel Q$$

**Erasure Function:**

$$|-| : \textsc{MonitoredSys} \rightarrow \textsc{Sys}$$

**Lemma**

$$M \longrightarrow_{\mathsf{mon}} M' \quad \text{implies} \quad |M| \longrightarrow |M'|$$

Partial Log Extraction

$$\textbf{pLog} : \kappa \rightarrow \mathbb{P}(\phi)$$

$$\textbf{pLog}(\epsilon) = \emptyset$$
$$\textbf{pLog}(\textbf{rcv}(p, \kappa); \kappa') = \textbf{rcv}(p); \textbf{pLogV}(\kappa') \cup \textbf{pLog}(\kappa)$$
$$\textbf{pLog}(\textbf{snd}(p, \kappa); \kappa') = \textbf{snd}(p); \textbf{pLogV}(\kappa') \cup \textbf{pLog}(\kappa)$$

$$\textbf{pLogV}(\epsilon) = \emptyset$$
$$\textbf{pLogV}(\textbf{rcv}(p, \kappa); \kappa') = \textbf{rcv}(p); \textbf{pLogV}(\kappa')$$
$$\textbf{pLogV}(\textbf{snd}(p, \kappa); \kappa') = \textbf{snd}(p); \textbf{pLogV}(\kappa')$$

**Definition**

$M$ has correct provenance iff $\forall \phi \in \textbf{pLog}(\textbf{prov}(M))$ we have $\phi \leq \textbf{log}(M)$.

**Theorem**

$M$ has correct provenance and $M \longrightarrow_{mon} M'$ implies $M'$ has correct provenance.

# Outline

# Conclusions

- Designed a provenance based calculus for distributed computing.
- Proposed a two-tier system for provenance tracking and usage.
- Defined provenance correctness
- Proved provenance correctness for our provenance tracking semantics.

# Conclusions

Thank You... Questions?