

Improving Data Quality for Web Services Composition

Xitong Li, Stuart Madnick, Hongwei Zhu, Yushun Fan

Stuart E. Madnick

MIT Sloan School of Management

smadnick@mit.edu



The 7th International Workshop on Quality in Databases (QDB'09)
August 24th, 2009, Lyon, France

Agenda

- **Research Motivation** – *focused on data quality*
 - Problem of **data misinterpretation**:
Different Web services have different assumptions about the interpretation of the exchanged data among them in the composition
- **Our Solution Approach**
 - Representation of Ontology and Context
 - Semantic Annotation
 - Reconciliation Approach
 - Prototype: *CMT*
 - Summary and Future Work

Real-world Scenarios:

GetQuote operation of the StockQuote service

The screenshot shows a web browser window with the URL <http://seekda.com/providers/webservicex.com/StockQuote?tab=usenow>. The page features a three-step process:

- Step 1:** Fill the input values.
- Step 2:** Sending Request to webservicex.com.
- Step 3:** View Result and leave comment.

The main section is titled "Web Service Tester" and includes a "privacy policy" link. It states: "You are invoking the **GetQuote** operation of the **StockQuote** service. Your request will be send to: <http://www.webservicex.com/stockquote.asmx>".

The operation being tested is "Get Stock quote for a company Symbol". The input field is labeled "symbol type:string" and contains the text "IBM".

At the bottom of the form, there are two buttons: "Go" and "Reset".

Response for IBM

Step 1 Fill the input values. Step 2 Sending Request to webservicex.com. Step 3 View Result and leave comment.

HTTP Response

```
<StockQuotes>
<Stock>
<Symbol>IBM</Symbol>
<Last>117.13</Last>
<Date>7/24/2009</Date>
<Time>9:44am</Time>
<Change>+0.07</Change>
<Open>116.68</Open>
<High>117.46</High>
<Low>116.63</Low>
Volume: 568498 /Volume:
<MktCap>153.4B</MktCap>
<PreviousClose>117.06</PreviousClose>
<PercentageChange>+0.06%</PercentageChange>
<AnnRange>69.50 - 130.93</AnnRange>
<Earnings>9.368</Earnings>
<P-E>12.50</P-E>
<Name>INTL BUSINESS MAC</Name>
</Stock>
</StockQuotes>
```

<MktCap> 153.4B

1. What currency it uses?
2. What does "B" mean?

(Even if you knew "B" meant Billion – is it US billion or UK billion?)

Share your opinion on this service: Cast your vote 1-5: ★★★★★ (cast your vote by clicking the corresponding star)

Tell other people what is your opinion:

View [all existing comments.](#)

Another Example: Response for ITWO

The screenshot shows a web browser window with the URL `http://seekda.com/providers/webservicex.com/StockQuote?tab=usenow#`. A modal window titled "HTTP Response" displays the following XML data:

```
<StockQuotes>
<Stock>
<Symbol>ITWO</Symbol>
<Last>13.18</Last>
<Date>7/24/2009</Date>
<Time>9:57 am</Time>
<Change>-0.20</Change>
<Open>13.10</Open>
<High>13.27</High>
<Low>13.09</Low>
<Volume>7052</Volume>
<MktCap>289.8M</MktCap>
<PreviousClose>13.38</PreviousClose>
<PercentageChange>-1.49%</PercentageChange>
<AnnRange>5.50 - 14.60</AnnRange>
<Earnings>3.973</Earnings>
<P-E>3.37</P-E>
<Name>i2 Technologies</Name>
</Stock>
</StockQuotes>
</string>
```

Red annotations on the XML include a circle around the symbol "ITWO" and a box around the market cap "<MktCap>289.8M". Red text asks:

- 1. What currency it uses?
- 2. What does "M" mean?

Different Example: Xignite Web Service

The screenshot shows a web browser window displaying the Xignite website. The page is titled "U.S. SEC EDGAR Filings Web Service - XigniteEdgar - Market Data Feeds, Financial Web Services :: Xignite - Maxthon 2.1.5". The URL in the address bar is "http://www.xignite.com/xEdgar.aspx?op=GetTotalAssets". The page features a navigation menu with links for Products, Solutions, Clients, Partners, News, Support, About, and Blog. The main content area is titled "XigniteEdgar" and "U.S. SEC EDGAR Filings". The operation being viewed is "GetTotalAssets", which is highlighted in red. Below the operation name, there are tabs for "Test Form", "Inputs", "Outputs", "SOAP", "GET", "POST", and "Sample Code". The "Description" section states: "Returns total asset information extracted from latest 10-Q or 10-K statements." The "Hits" section indicates that requests against this operation count as one hit. The "Test Form" section provides instructions on how to use the form and lists three bullet points: "Sample values are filled out for convenience.", "Results are displayed in a new window.", and "You must sign up for Free Trial to use this form." Below the test form, there is a table with two columns: "Parameter" and "Value". The "Identifier" parameter is set to "MSFT" and the "IdentifierType" parameter is set to "Symbol". At the bottom of the page, there are buttons for "View Results As:" with options for "XML", "Tabular Format", and "Spreadsheet". The right sidebar contains a search bar and a list of operations including "Overview", "WSDL", "Pricing", "FAQ", "Demo", "Operations", "ListSICCodes", "ListSECFilingTypes", "LookupCIK", "GetCIK", "SearchFilings", "GetFilingOccurrences", "PredefinedQueryFilings...", "QueryFilings", "QueryMasterDocuments", and "GetLastFiling".

U.S. SEC EDGAR Filings Web Service - XigniteEdgar - Market Data Feeds, Financial Web Services :: Xignite - Maxthon 2.1.5

http://www.xignite.com/xEdgar.aspx?op=GetTotalAssets

My Cart My Account Logout
TOLL FREE: 1-888-XML-SOAP
(1-888-965-7827)

Products Solutions Clients Partners News Support About Blog

Home > Products > Company Data On-Demand > XigniteEdgar

XigniteEdgar

U.S. SEC EDGAR Filings

Operation: **GetTotalAssets**

Test Form Inputs Outputs SOAP GET POST Sample Code

Description

Returns total asset information extracted from latest 10-Q or 10-K statements.

Hits

Requests against this operation count as one hit.

Test Form

This form lets you test the operation using your browser:

- Sample values are filled out for convenience.
- Results are displayed in a new window.
- You must sign up for **Free Trial** to use this form.

Parameter	Value
Identifier	MSFT
IdentifierType	Symbol

View Results As: XML Tabular Format Spreadsheet

search the catalog

free trial

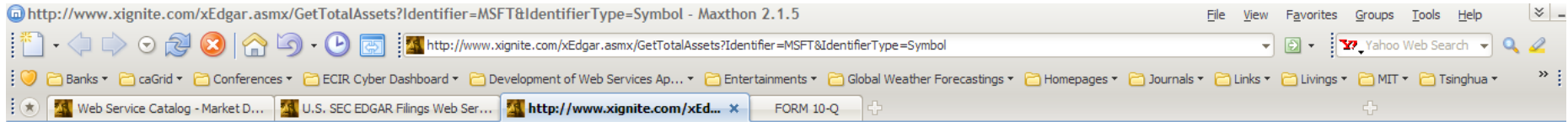
buy online

Overview
WSDL
Pricing
FAQ
Demo
Operations
ListSICCodes
ListSECFilingTypes
LookupCIK
GetCIK
SearchFilings
GetFilingOccurrences
PredefinedQueryFilings...
QueryFilings
QueryMasterDocuments
GetLastFiling

In the 1st year of marriage, man speaks and woman listens. In the 2nd year, woman speaks and man listens. In the 3rd year, they both speak and neighbors listen.

1674M Zoom: 100%

GetTotalAssets for MSFT



```
<?xml version="1.0" encoding="utf-8" ?>
- <TotalAssets xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema" xmlns="http://www.xignite.com/services/">
  <Outcome>Success</Outcome>
  <Identity>Cookie</Identity>
  <Delay>0.016</Delay>
  - <Security>
    <Outcome>Success</Outcome>
    <Delay>0</Delay>
    <CIK>0000789019</CIK>
    <Cusip>594918104</Cusip>
    <Symbol>MSFT</Symbol>
    <ISIN>US5949181045</ISIN>
    <Valoren>951692</Valoren>
    <Name>Microsoft Corporation</Name>
    <Market>NASDAQGS</Market>
    <CategoryOrIndustry>TECHNOLOGY</CategoryOrIndustry>
  </Security>
  <Source>10-Q/K</Source>
  <SourceDate>04/23/2009</SourceDate>
  <SourceUrl>http://www.sec.gov/Archives/edgar/data/789019/000119312509085779/d10q.htm</SourceUrl>
  <SourceType>Text</SourceType>
  <Value>68853</Value>
</TotalAssets>
```

What is this date 04/23/2009 ?
What if it was 04/05/06 ?

MSFT Total Assets: 68,853

GetTotalAssets for ITWO

```
http://www.xignite.com/xEdgar.aspx/GetTotalAssets?Identifier=ITWO&IdentifierType=Symbol - Maxthon 2.1.5
http://www.xignite.com/xEdgar.aspx/GetTotalAssets?Identifier=ITWO&IdentifierType=Symbol
<?xml version="1.0" encoding="utf-8" ?>
- <TotalAssets xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema" xmlns="http://www.xignite.com/services/">
  <Outcome>Success</Outcome>
  <Identity>Cookie</Identity>
  <Delay>0.031</Delay>
- <Security>
  <Outcome>Success</Outcome>
  <Delay>0</Delay>
  <CIK>0001009304</CIK>
  <Cusip>465754208</Cusip>
  <Symbol>ITWO</Symbol>
  <ISIN>US4657542084</ISIN>
  <Valoren>2074416</Valoren>
  <Name>i2 Technologies, Inc.</Name>
  <Market>NASDAQGM</Market>
  <CategoryOrIndustry>TECHNOLOGY</CategoryOrIndustry>
</Security>
  <Source>10-Q/K</Source>
  <SourceDate>05/07/2009</SourceDate>
  <SourceUrl>http://www.sec.gov/Archives/edgar/data/1009304/000119312509103105/d10q.htm</SourceUrl>
  <SourceType>Text</SourceType>
  <Value>313776</Value>
</TotalAssets>
```

MSFT Total Assets: 68,853 vs. ITWO Total Assets: 313,776

?

(ITWO has Five Times more assets than MSFT?)

[Table of Contents](#)

MICROSOFT CORPORATION

BALANCE SHEETS
(In millions)

(in millions)

	March 31, 2009 (Unaudited)	June 30, 2008(1)
Assets		
Current assets:		
Cash and cash equivalents	\$ 7,285	\$ 10,339
Short-term investments (including securities pledged as collateral of \$1,445 and \$2,491)	18,055	13,323
Total cash, cash equivalents, and short-term investments	25,340	23,662
Accounts receivable, net of allowance for doubtful accounts of \$242 and \$153	9,182	13,589
Inventories	657	985
Deferred income taxes	1,926	2,017
Other	3,619	2,989
Total current assets	40,724	43,242
Property and equipment, net of accumulated depreciation of \$7,236 and \$6,302	7,112	6,242
Equity and other investments	4,112	6,588
Goodwill	12,554	12,108
Intangible assets, net	1,756	1,973
Deferred income taxes	956	949
Other long-term assets	1,039	1,691
Total assets	\$ 68,853	\$ 72,793
Liabilities and stockholders' equity		
Current liabilities:		
Accounts payable	\$ 3,017	\$ 4,034
Short-term debt	1,999	-
Accrued compensation	2,644	2,934
Income taxes	773	3,248
Short-term unearned revenue	10,924	13,397
Securities lending payable	1,533	2,614
Other	2,933	3,659
Total current liabilities	23,823	29,886
Long-term unearned revenue	1,388	1,900
Other long-term liabilities	6,699	4,721
Commitments and contingencies		
Stockholders' equity:		
Common stock and paid-in capital – shares authorized 24,000; outstanding 8,898 and 9,151	61,896	62,849
Retained deficit, including accumulated other comprehensive income of \$726 and \$1,110	(7,039)	(1,056)

[Table of Contents](#)

PART 1. FINANCIAL INFORMATION

ITEM 1. FINANCIAL STATEMENTS

i2 TECHNOLOGIES, INC.
CONDENSED CONSOLIDATED BALANCE SHEETS
 (In thousands, except par value)
 (unaudited)

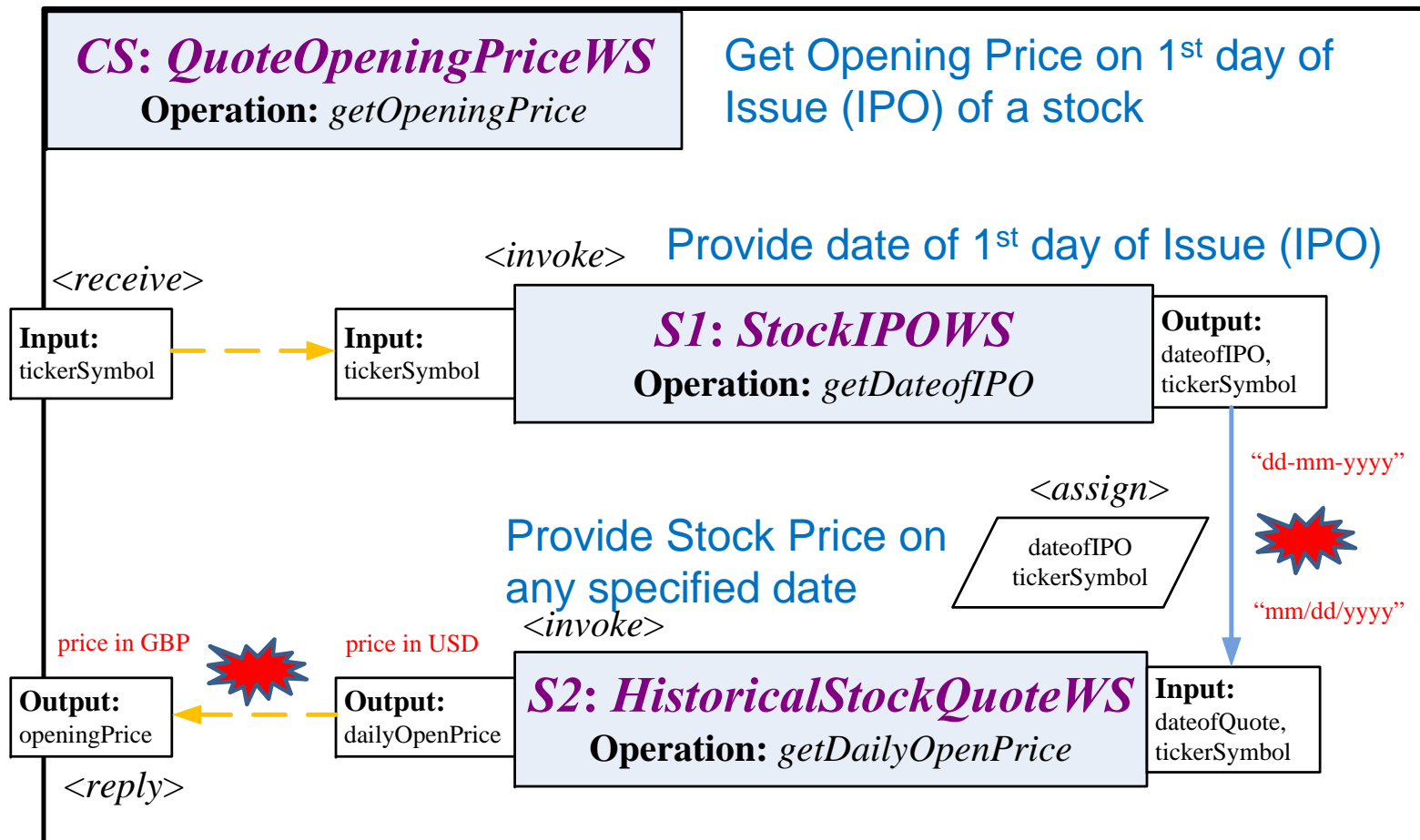
(in thousands)

	March 31, 2009	December 31, 2008 <small>(as restated, see Note 10)</small>
ASSETS		
Current assets:		
Cash and cash equivalents	\$ 157,659	\$ 238,013
Restricted cash	8,921	5,777
Accounts receivable, net	22,980	25,846
Other current assets	7,418	9,477
Total current assets	196,978	279,113
Premises and equipment, net	4,035	4,915
Goodwill	16,684	16,684
Non-current deferred tax asset	5,846	7,289
Other non-current assets	3,371	5,024
Total assets	\$ 226,914	\$ 313,025
LIABILITIES AND STOCKHOLDERS' EQUITY		
Current liabilities:		
Accounts payable	\$ 3,878	\$ 4,855
Accrued liabilities	16,499	15,116
Accrued compensation and related expenses	11,489	18,679
Deferred revenue	58,597	53,028
Total current liabilities	90,463	91,678
Total long-term debt, net	—	64,520
Taxes payable	5,292	6,948
Total liabilities	95,755	163,146
Commitments and contingencies		
Stockholders' equity:		
Preferred Stock, \$0.001 par value, 5,000 shares authorized, none issued and outstanding	—	—
Series A junior participating preferred stock, \$0.001 par value, 2,000 shares authorized, none issued and outstanding	—	—

Research Motivation

- **Goal of Using Web Services**
 - Web services composition/integration/mashup
 - **Combine multiple web services to create a composite web service**
 - Standards: WSDL (single WS), BPEL (composition process)
- **Long-Standing Challenge for Data Quality**
 - Problems of **data misinterpretation**
 - Different Web services have different assumptions about the interpretation of the exchanged data among them in the composition
 - Root cause: inconsistent data representation, unit, precision, scaling, and meaning, etc.

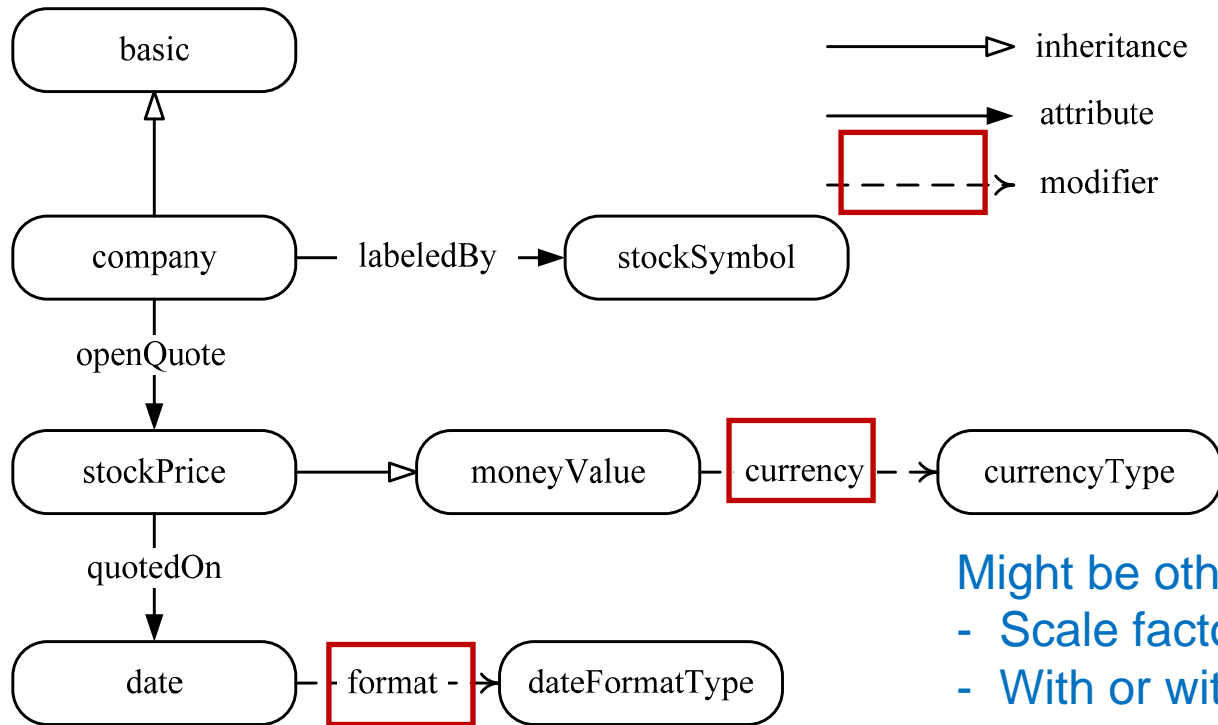
Example Used in our Demonstration: Illustrating Creation of Composite Web Service



- Notes: 1. Assumptions of data interpretation are not explicitly represented in WSDLs.
2. Often parallel asynchronous processing (omitted from example)

Context-enriched Ontology Model

- **Ontology with Context Modifiers**



Might be other modifiers, e.g.:

- Scale factor
- With or without fees

Service	Context Definition
CS, S1	$ctxtUK = \{ \langle format, dd-mm-yyyy \rangle, \langle currency, GBP \rangle \}$
S2	$ctxtUS = \{ \langle format, mm/dd/yyyy \rangle, \langle currency, USD \rangle \}$

Specification of Conversion Among Modifier Types (generic conversions – not specific to example)

- **Conversion Representation**

$cvt(C, m, ctxt_s, ctxt_t, mvs, mvt, vs, vt)$

$cvt_{format}(date, format, ctxtUK, ctxtUS, \underline{\text{"dd-mm-yyyy"}}, \underline{\text{"mm/dd/yyyy"}}, \text{"09-08-2009"}, vt)$
 $vt = \text{"08/09/2009"}$

- **Conversion Implementation**

- Xpath functions: cvt_{format}

$Vt = \text{Concat}(\text{substring-before}(\text{substring-after}(Vs, \text{"-"}), \text{"-"}), \text{"/"},$
 $\text{substring-before}(Vs, \text{"-"}), \text{"/"}, \text{substring-after}(\text{substring-after}(Vs, \text{"-"}), \text{"-"}))$

- External services: $cvt_{currency}$

- CurrencyExchangeService

Semantic Annotation of WSDL (using SAWSDL)

```
<wsdl:definitions targetNamespace="http://stockQuote.coin.mit"
  xmlns:stkOntology="http://stockQuote.coin.mit/ontologies/stockOntology#" xmlns:sawSDL="http://www.w3.org/ns/sawSDL" ...
  sawSDL:modelReference="stkOntology#ctxtUK">
  <wsdl:types>
  <schema elementFormDefault="qualified" targetNamespace="http://stockQuote.coin.mit" xmlns="http://www.w3.org/2001/XMLSchema">
    <element name="getDateofIPO">
      <complexType>
        <sequence>
          <element name="tickerSymbol" type="xsd:string" sawSDL:modelReference="stkOntology#stockSymbol"/>
        </sequence>
      </complexType>
    </element>
    <element name="getDateofIPOResponse">
      <complexType>
        <sequence>
          <element name="getDateofIPOReturn" type="impl:IPOBean"/>
        </sequence>
      </complexType>
    </element>
    <complexType name="IPOBean">
      <sequence>
        <element name="dateofIPO" nillable="true" type="xsd:string" sawSDL:modelReference="stkOntology#date stkOntology#ctxtUK"/>
        <element name="tickerSymbol" nillable="true" type="xsd:string" sawSDL:modelReference="stkOntology#stockSymbol"/>
      </sequence>
    </complexType>
  </schema>
</wsdl:types>
```

Both **global** and **local** context specification possible.

Reconciliation Approach for Any Web Service Composition

*Generic Ontology/SAWSDL context definitions of
web services and Conversion Specifications set up*

Our Reconciliation Approach

1. Translating WSDL/BPEL to LOTOS NT
2. Detecting Context Conflicts
3. Incorporating Conversions into LOTOS NT
4. Generating Mediated BPEL

1. Translating WSDL/BPEL to LOTOS NT

- **Benefits of LOTOS NT**

- Formalism for verifying composition processes, e.g., deadlock-freeness
- Independence of any composition languages, e.g., BPEL, OWL-S process

- **Translating Static Aspect to Types**

- WSDL: data types, messages, operations
- BPEL: variables, partner links

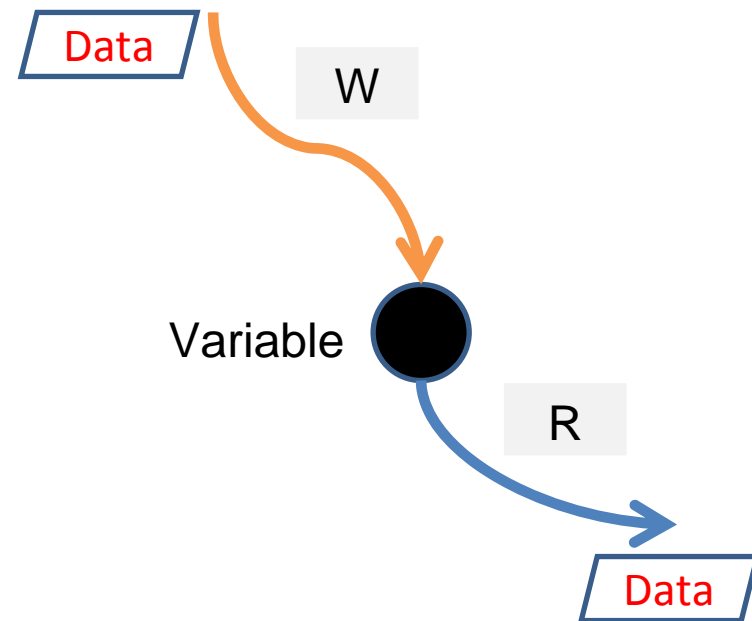
```
(* @ tickerSymbol:mdlRef="stkOntology#stockSymbol"  
* @ dateofIPO:mdlRef="stkOntology#dateofIPO" stkOntology#ctxtUK"  
*)  
type IPOBeanComplexType is  
  IPOBeanComplexType (tickerSymbol:string, dateofIPO:string)  
end type
```

2. Detecting Context Conflicts

- 2.1 Identify Data Transfers

- Explicit data transfers: `<assign>` (BPEL), “:=” (LOTOS NT)
- Implicit data transfers through variables in the process

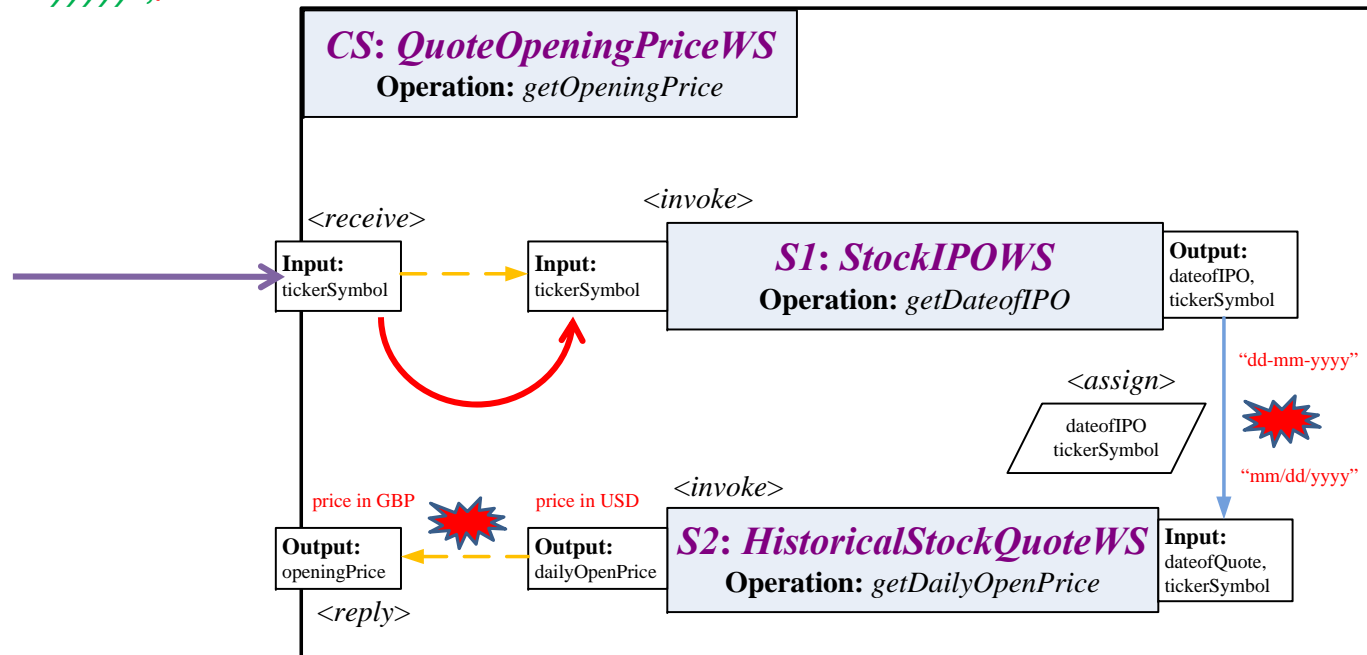
Implicit data transfers refer to the situations that *a piece of data is written to a variable and then directly read from the variable, without assignments between the write and read operations.*



2. Detecting Context Conflicts

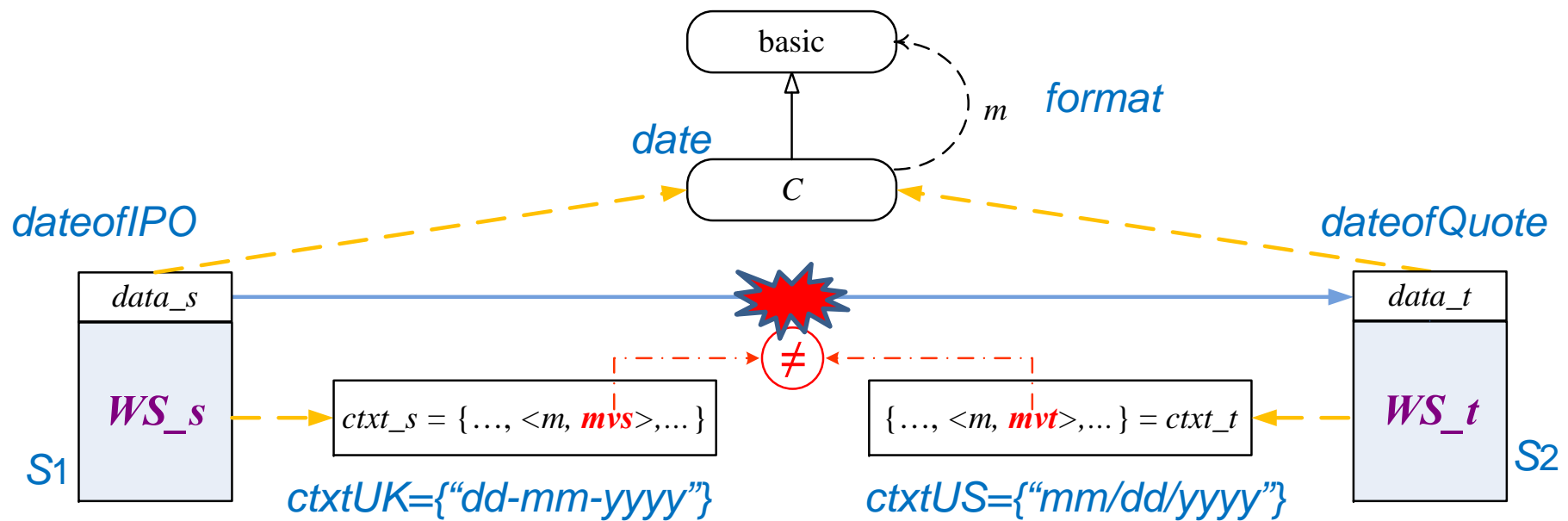
- **Implicit Data Transfers**

```
QuoteOpeningPricePL(?quoteOpeningPriceRole
(getOpeningPriceInput(getOpeningPriceRequestMessage
(getDateofIPO)))) ;
StockIPOPL(!stockIPOServiceRole(getDateofIPOInput
(getDateofIPORequestMessage(getDateofIPOType
(getDateofIPO))))
```



2. Detecting Context Conflicts

- 2.2 Examine Each Data Transfer for Conflicts



Typical Scenario of Context Conflicts within a Data Transfer

2. Detecting Context Conflicts

- Detected Context Conflicts

	Source	Target
LOTOS NT	dateofIPO	dateofQuote
Web service	S1	S2
Context	ctxtUK	ctxtUS
date.format	“dd-mm-yyyy”	“mm/dd/yyyy”

	Source	Target
LOTOS NT	getDailyOpenPriceResponse	getDailyOpenPriceResponse
Web service	S2	CS
Context	ctxtUS	ctxtUK
moneyValue.currenc	“USD”	“GBP”

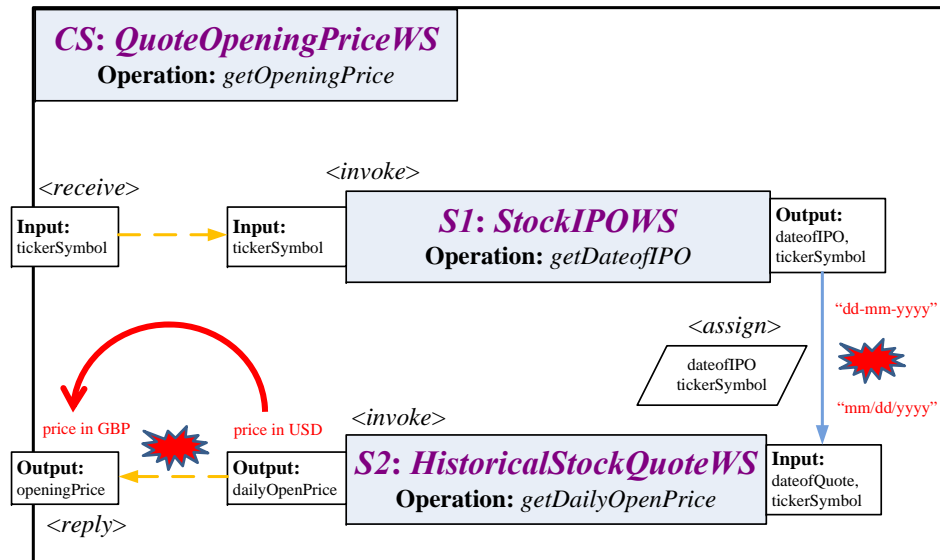
y

3. Incorporating Conversions into LOTOS NT

3.1 Make implicit data transfers with conflicts explicit


```

HistoricalStockQuotePL (?historicalStockQuoteServiceRole
(getDailyOpenPriceOutput(getDailyOpenPriceResponseMessage
(getDailyOpenPriceResponse)))) ;
QuoteOpeningPricePL(!quoteOpeningPriceRole
(getOpeningPriceOutput(getOpeningPriceResponseMessage
(getDailyOpenPriceResponse))))
    
```



3. Incorporating Conversions into LOTOS NT

```
HistoricalStockQuotePL (?historicalStockQuoteServiceRole  
(getDailyOpenPriceOutput(getDailyOpenPriceResponseMessage  
(getDailyOpenPriceResponse)))) ;  
QuoteOpeningPricePL(!quoteOpeningPriceRole  
(getOpeningPriceOutput(getOpeningPriceResponseMessage  
(getDailyOpenPriceResponse))))
```

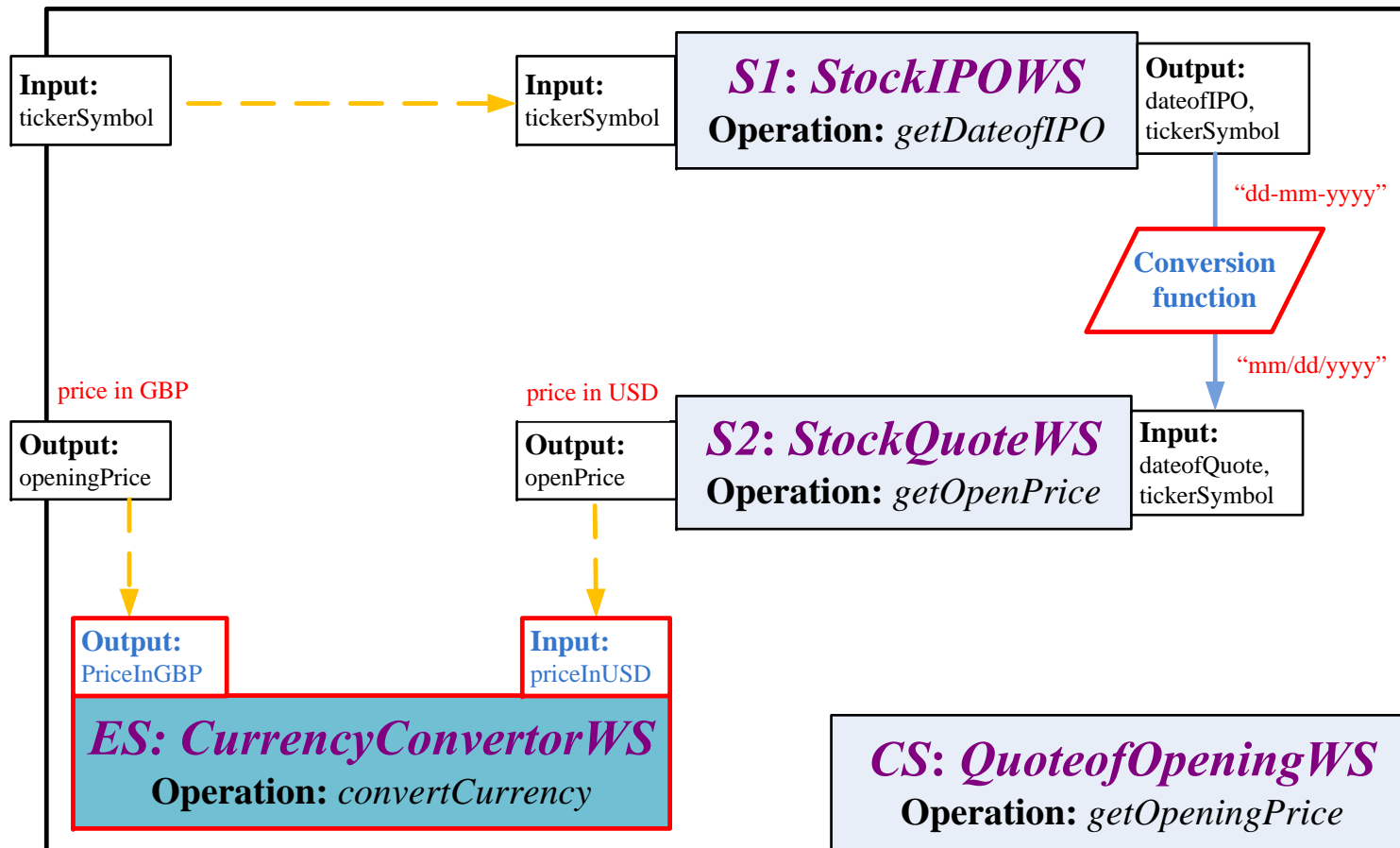


```
var getOpeningPriceResponse
```

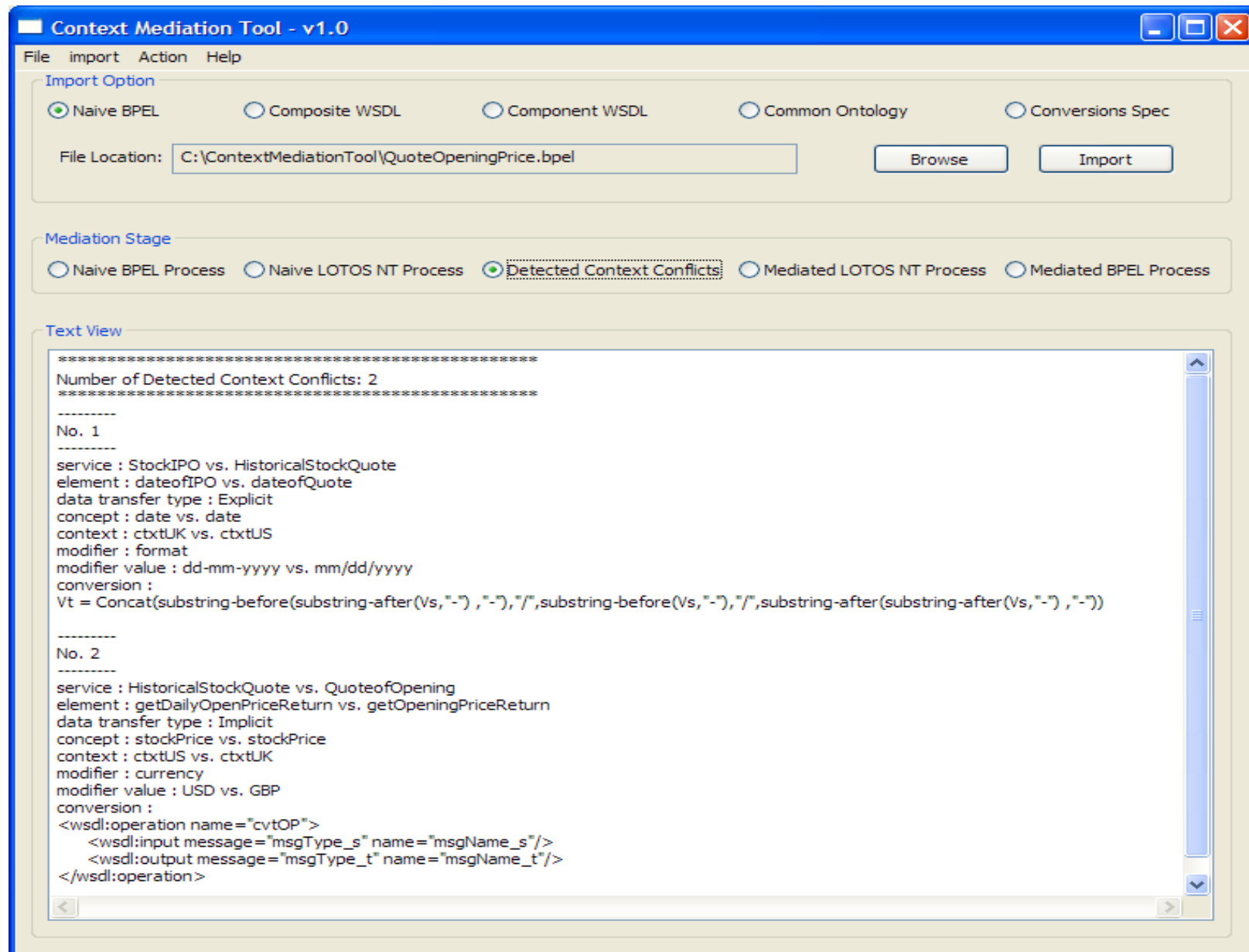
```
HistoricalStockQuotePL (?historicalStockQuoteServiceRole  
(getDailyOpenPriceOutput(getDailyOpenPriceResponseMessage  
(getDailyOpenPriceResponse)))) ;  
getOpeningPriceResponse := getDailyOpenPriceResponse;  
QuoteOpeningPricePL(!quoteOpeningPriceRole  
(getOpeningPriceOutput(getOpeningPriceResponseMessage  
(getOpeningPriceResponse ))))
```

3. Incorporating Conversions into LOTOS NT

- 3.2 Incorporate conversions and generate mediated process



Prototype: CMT – Illustrating Conflict Detection



Summary and Future Work

- **Strengths**

- Address data misinterpretation problems to improve quality
- Separation of concern between service descriptions (e.g., BPEL) and ontology/context descriptions
- Use the standard-compliant method (i.e., SAWSDL) for semantic annotation
- Use LOTOS, a kind of process algebra, which can be integrated with behavioral compatibility verification
- Automates the incorporation of conversions into composition (can assemble complex conversions)

- **Ongoing Work**

- Extend to deal with semantic heterogeneity at both schematic and data instance level
- (Semi-)automatic construction of ontology, contexts and annotations

Maybe Web Services were in use during the Napoleonic Wars? *(Data Quality Consequences)*

In 1805, the Austrian and Russian Emperors agreed to join forces against Napoleon. The Russians promised that their forces would be in the field in Bavaria by **Oct. 20**.

The Austrian staff planned its campaign based on that date in the **Gregorian calendar**. Russia, however, still used the ancient **Julian calendar**, which lagged 10 days behind.

The calendar difference allowed Napoleon to surround Austrian General Mack's army at Ulm and force its surrender on Oct. 21, well before the Russian forces could reach him, ultimately setting the stage for Austerlitz.

Source: David Chandler, *The Campaigns of Napoleon*, New York: MacMillan 1966, pg. 390.

**Thanks again for your interest in this work.
Please refer to the paper for more details.**

Any questions?

Backup Slides

1. Translating WSDL/BPEL to LOTOS NT

- Translating Dynamic Aspect to Process

BPEL	LOTOS NT
<i><receive variable="v" .../></i>	<i>(c?v)</i>
<i><reply variable="v" .../></i>	<i>(c!v)</i>
<i><invoke inputVariable="v1" outputVariable="v2" .../></i>	<i>(c!v1 ; c?v2)</i>
<i><assign ...> <copy> <from variable="v1"> <to variable="v2"/> </copy> <copy> <from variable="v3"> <to variable="v4"/> </copy> </assign></i>	<i>(v2 := v1; v4 := v3)</i>
<i><sequence ...> <...activity1.../> <...activity2.../> </sequence></i>	<i>(action1 ; action2)</i>

1. Translating WSDL/BPEL to LOTOS NT

```

QuoteOpeningPricePL(?quoteOpeningPriceRole
(getOpeningPriceInput(getOpeningPriceRequestMessage
(getDateofIPO)))) ;
StockIPOPL(!stockIPOServiceRole(getDateofIPOInput
(getDateofIPORequestMessage(getDateofIPOType
(getDateofIPO)))) ;
StockIPOPL(?stockIPOServiceRole(getDateofIPOOutput
(getDateofIPOResponseMessage(getDateofIPOResponse)))) ;
getDailyOpenPrice := getDailyOpenPriceType (getDateofIPOResponse.getDateofIPOReturn.tickerSymbol,
getDateofIPOResponse.getDateofIPOReturn.dateofIPO) ;
HistoricalStockQuotePL(!historicalStockQuoteServiceRole
(getDailyOpenPriceInput(getDailyOpenPriceRequestMessage (getDailyOpenPrice)))) ;
HistoricalStockQuotePL (?historicalStockQuoteServiceRole
(getDailyOpenPriceOutput(getDailyOpenPriceResponseMessage (getDailyOpenPriceResponse)))) ;
QuoteOpeningPricePL(!quoteOpeningPriceRole
(getOpeningPriceOutput(getOpeningPriceResponseMessage (getDailyOpenPriceResponse))))
    
```

