



Wendelin Big Data

Industrial Monitoring Platform

2014-04-03 – Paris

Who are we?

- **Jean-Paul Smets**
- **Nexedi CEO**
- **Author of ERP5**
- **jp@nexedi.com**

- **Ivan Tyagov**
- **Senior Developer**
- **Wendelin project lead**
- **ivan@nexedi.com**

Who is missing?

- **Kirill Smelkov**
- **Senior Developer**
- **wendelin.core**
- **Sebastien Robin**
- **Project Director**
- **Author of POC**

Agenda


- **Where do we come from**
- **Wendelin Architecture**
- **Detailed Example**
- **Future Roadmap**

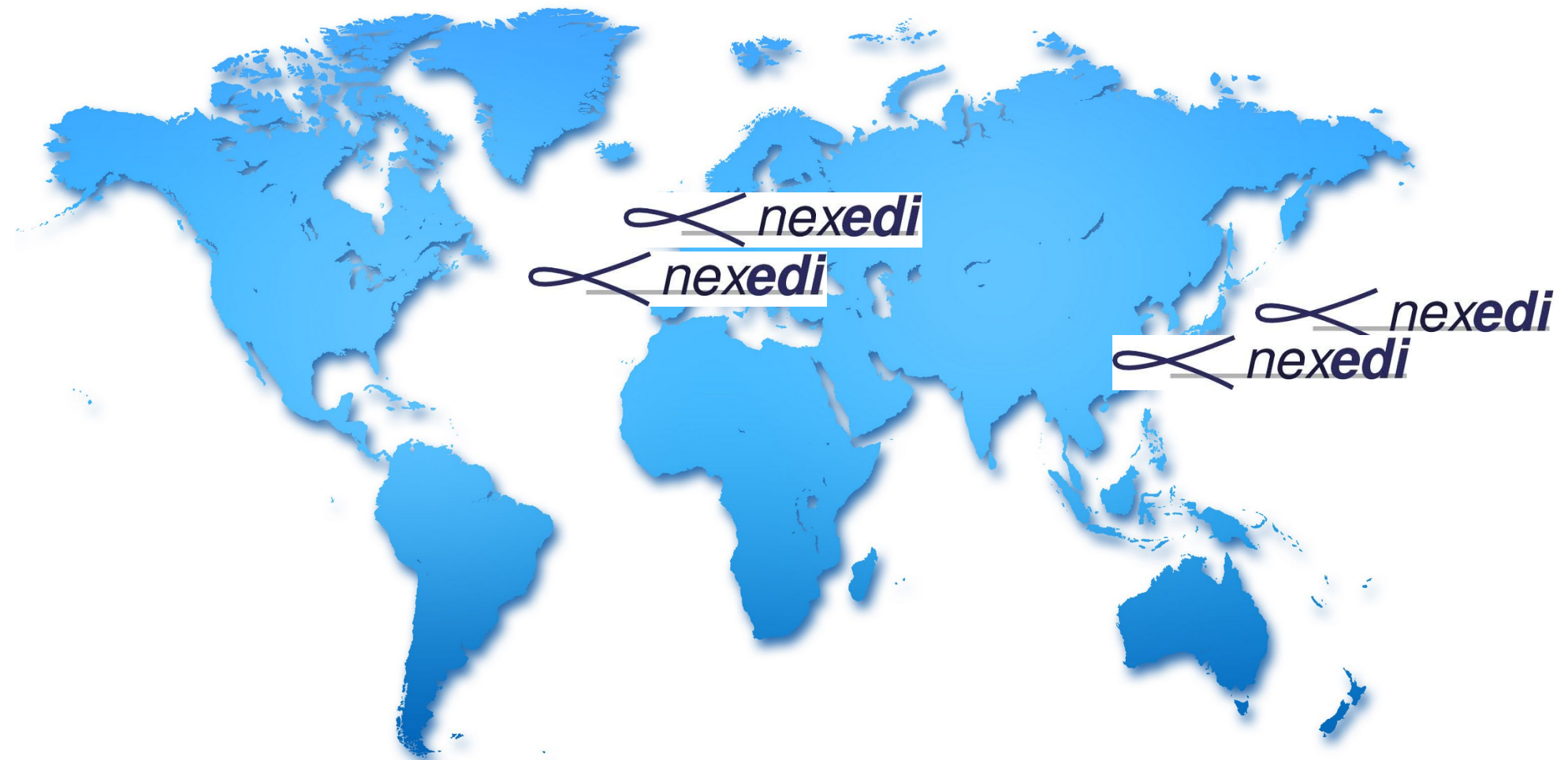
Where do we come from?



Nexedi

- **Possibly Largest OSS Publisher in Europe**

- ERP5: ERP, CRM, ECM, e-business framework
- SlapOS: distributed mesh cloud operation system
- NEO: distributed transactional NoSQL database
-  – **Wendelin: out-of-core big data based on NumPy**
- re6st: resilient IPv6 mesh overlay network
- RenderJS: javascript component system
- JIO: javascript virtual database and virtual filesystem
- cloudooo: multimedia conversion server
- Web Runner: web based Platform-as-a-Service (PaaS) and IDE
- OfficeJS: web office suite based on RenderJS and JIO





Application Convergence



+



?

Case 1: Wind Turbines



- **Collect logs**
- **Collect records**
- **Predict failure**
- **Plan maintenance**
- **Reduce downtime**
→ **add X% profits**

ERP5

scikit
learn

ERP5

Case 2: Cars



- **Collect logs**
- **Collect records**
- **Predict failure**
- **Plan maintenance**
- **Reduce downtime**
→ **increase loyalty**

ERP5



ERP5

Case 3: Solar Energy



- **Collect logs**
- **Collect records**
- **Predict degradation**
- **Plan maintenance**
- **Increase efficiency**
 - add X% to profits



Wendelin Architecture



Standard Hardware no router / no SAN



inspur 浪潮 | lenovo 联想 | CORETO

x 160

- 2 x 10 Gbps
- 2 x 6 core Xeon CPU
- 512 GB RAM
- 4 x 1 TB SSD
- 1 x M2090 GPU

+



x 32

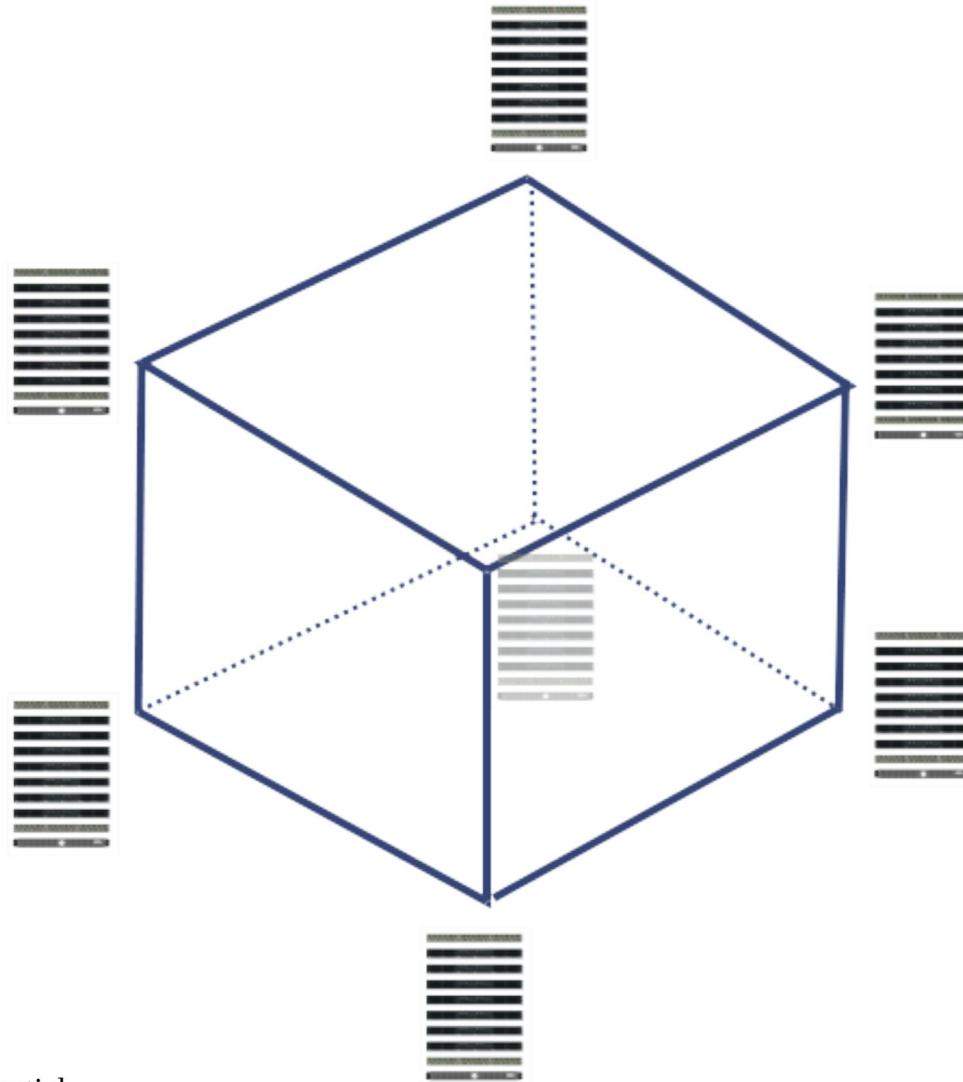
- 10 Gbps
- Unmanaged

+

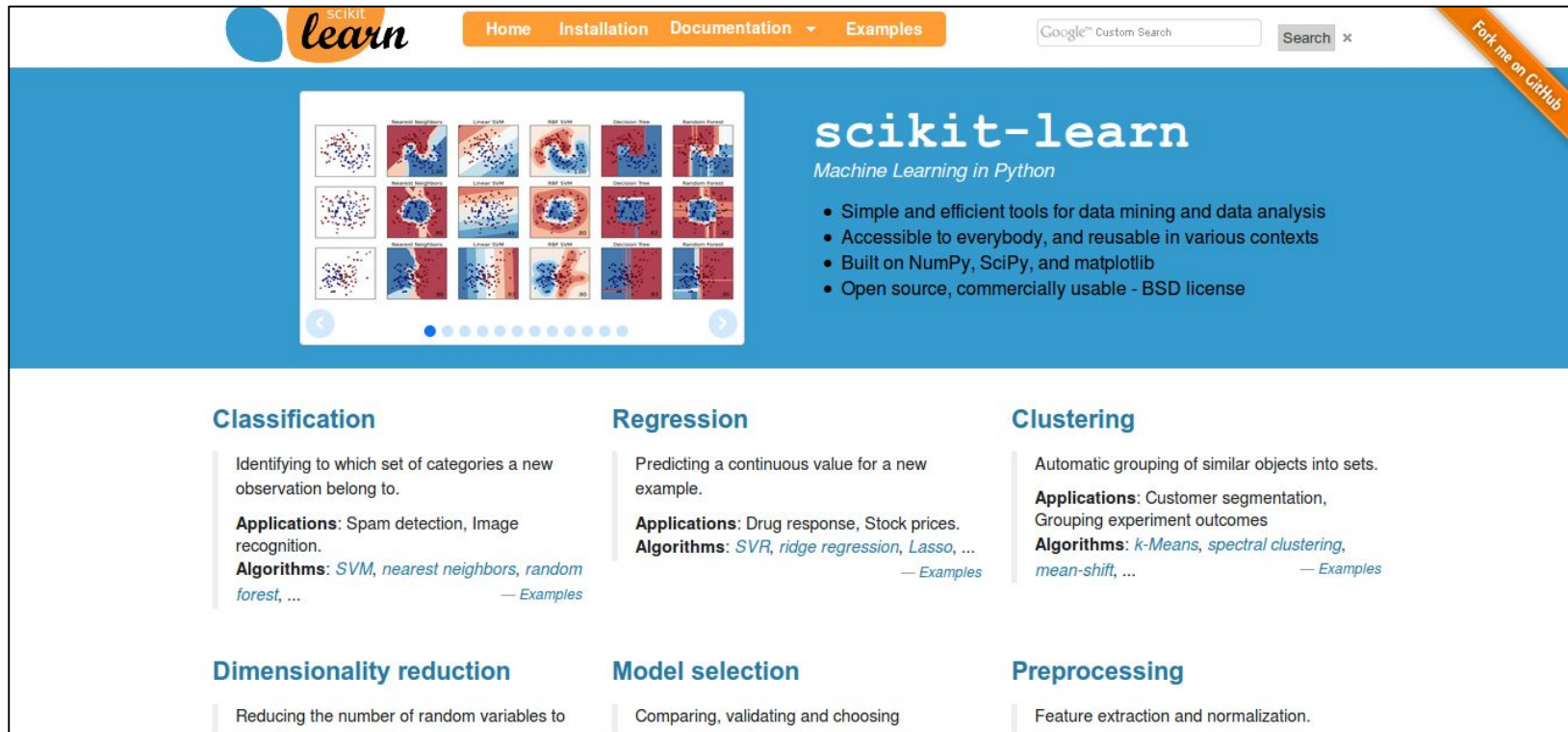


x 320

Wendelin Hypercube Datacenter



Take the Best Analytics scikit-learn.org



The screenshot shows the scikit-learn website homepage. At the top, there's a navigation bar with links for Home, Installation, Documentation, and Examples. A search bar is also present. The main header features the scikit-learn logo and the tagline "Machine Learning in Python". Below this, a grid of 12 small plots illustrates various machine learning models. To the right, a list of bullet points highlights the library's features: simple and efficient tools for data mining and data analysis, accessibility to everybody, reusability in various contexts, being built on NumPy, SciPy, and matplotlib, and being open source with a BSD license. A diagonal banner on the right side says "Fork me on GitHub".

scikit-learn
Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification
Identifying to which set of categories a new observation belong to.
Applications: Spam detection, Image recognition.
Algorithms: SVM, nearest neighbors, random forest, ... — Examples

Regression
Predicting a continuous value for a new example.
Applications: Drug response, Stock prices.
Algorithms: SVR, ridge regression, Lasso, ... — Examples

Clustering
Automatic grouping of similar objects into sets.
Applications: Customer segmentation, Grouping experiment outcomes
Algorithms: k-Means, spectral clustering, mean-shift, ... — Examples

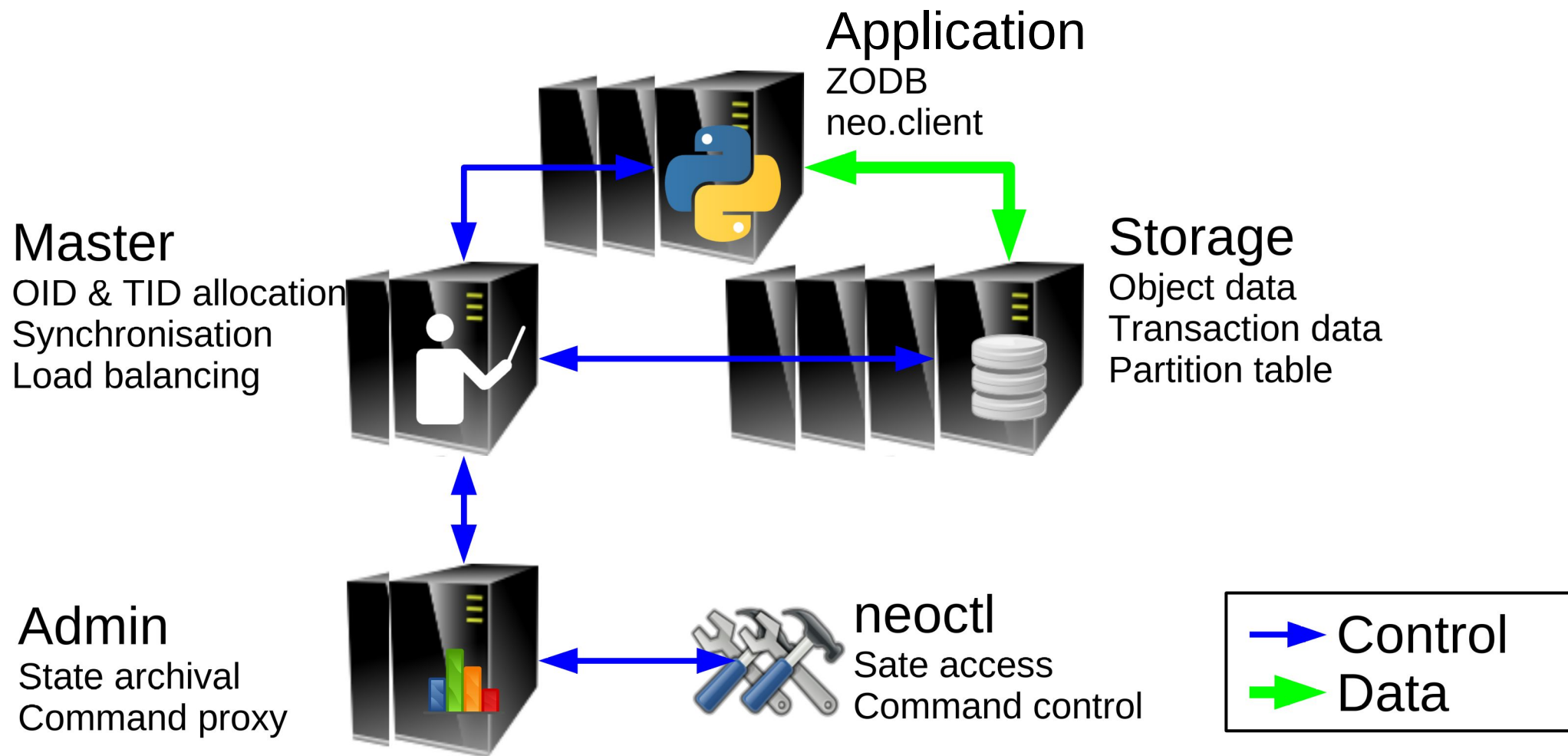
Dimensionality reduction
Reducing the number of random variables to

Model selection
Comparing, validating and choosing

Preprocessing
Feature extraction and normalization.



Add Distributed Storage neoppod.org



“Magic” out-of-core for NumPy

PyData Paris 2015 – 16h45 Kirill Smelkov

ZBigArray

1	2	3	4	5	6	7	8	9	10	11	12
---	---	---	---	---	---	---	---	---	----	----	----



1
5
9



2
6
10



3
7
11



4
8
12



Add Elastic PaaS erp5.com

```
# Initialize data
data_size = 1000000
server_count = 1000
chunk_size = data_size / server_count
data = array(data_size)

# Process data in parallel on each server (Map Reduce, Batch, etc.)
for server in server_count:
    data.activate().process(server*chunk_size, chunk_size)
```



And Multicloud Deployment slapos.org



```
0 [buildout]
1
2 extends =
3 # "slapos" stack describes basic things needed for 99.9% of SlapOS Software
4 # Releases
5 ../../stack/slapos.cfg
6 # Extend here component profiles, like openssl, apache, mariadb, curl...
7 # Or/and extend a stack (lamp, tomcat) that does most of the work for you
8 # In this example we only need the dash binary to run a simple "hello world"
9 # shell script.
10 ../../component/dash/buildout.cfg
11
12 parts =
13 # Call installation of slapos.cookbook egg defined in stack/slapos.cfg (needs
14 # in 99.9% of Slapos Software Releases)
15 slapos-cookbook
16 # Call creation of instance.cfg file that will be called for deployment of
17 # instance
18 template
19
20 # Download instance.cfg.in (buildout profile used to deployment of instance),
21 # replace all ${foo:bar} parameters by real values, and change ${foo:bar} to
22 # ${foo:bar}
23 [template]
24 recipe = slapos.recipe.template
25 url = ${:_profile_base_location_}/instance.cfg.in
26 output = ${buildout:directory}/instance.cfg
27 # MD5 checksum can be skipped for development (easier to develop), but must be filled for production
28 md5sum = 1fc461c00e86485bee77a942f39e3c43
29 mode = 0644
30
```

Save



MMC Rus



L'Education change le monde



Wendelin Platform 100% open source

100% Python

Scikit Learn

Data Analytics

NEO

Distributed Storage



ERP5

Elastic PaaS

SlapOS

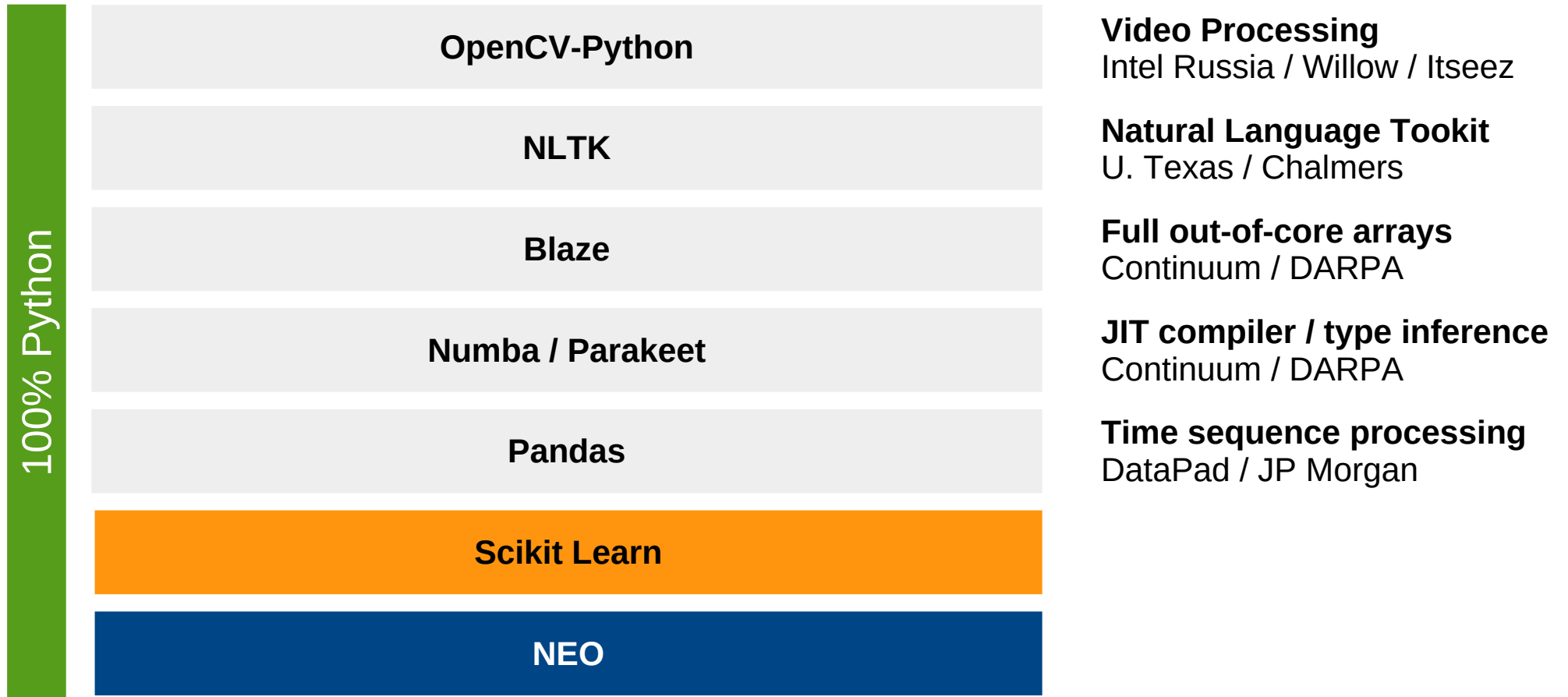
Multicloud Deployment



Multi Data Center

Wendelin Options

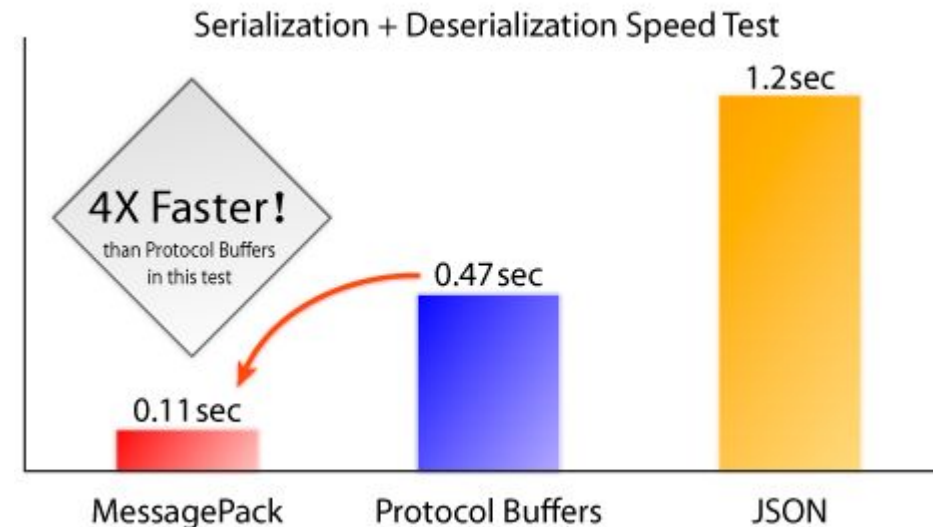
100% open source



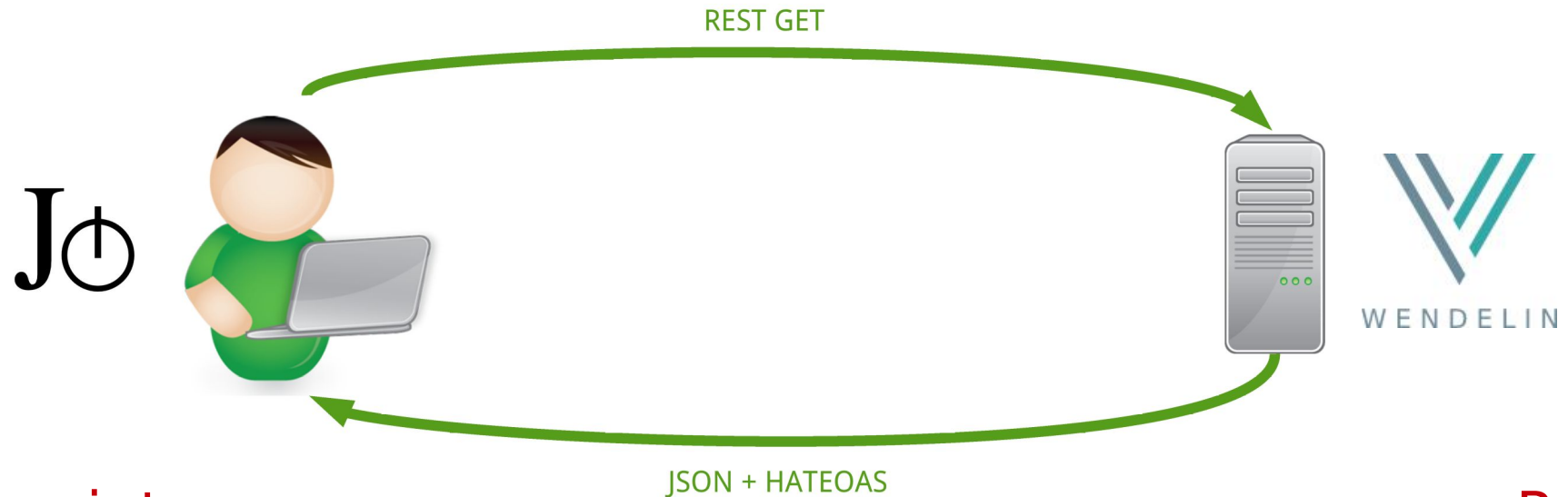
Data Ingestion: fluentd



- **Based on MsgPack middleware**
- **Created by TreasureData (BDaaS pioneers)**
- **Used by Amazon**
- **Numerous plugins**
- **Scalable and resilient**
- **Bandwidth saver**



Wendelin UI



Javascript

- **HTML5 Render** RenderJS
- **Data vizualisation**
- **Offline support** JIO

Python

- **Data access** REST API
- **Batch processing**

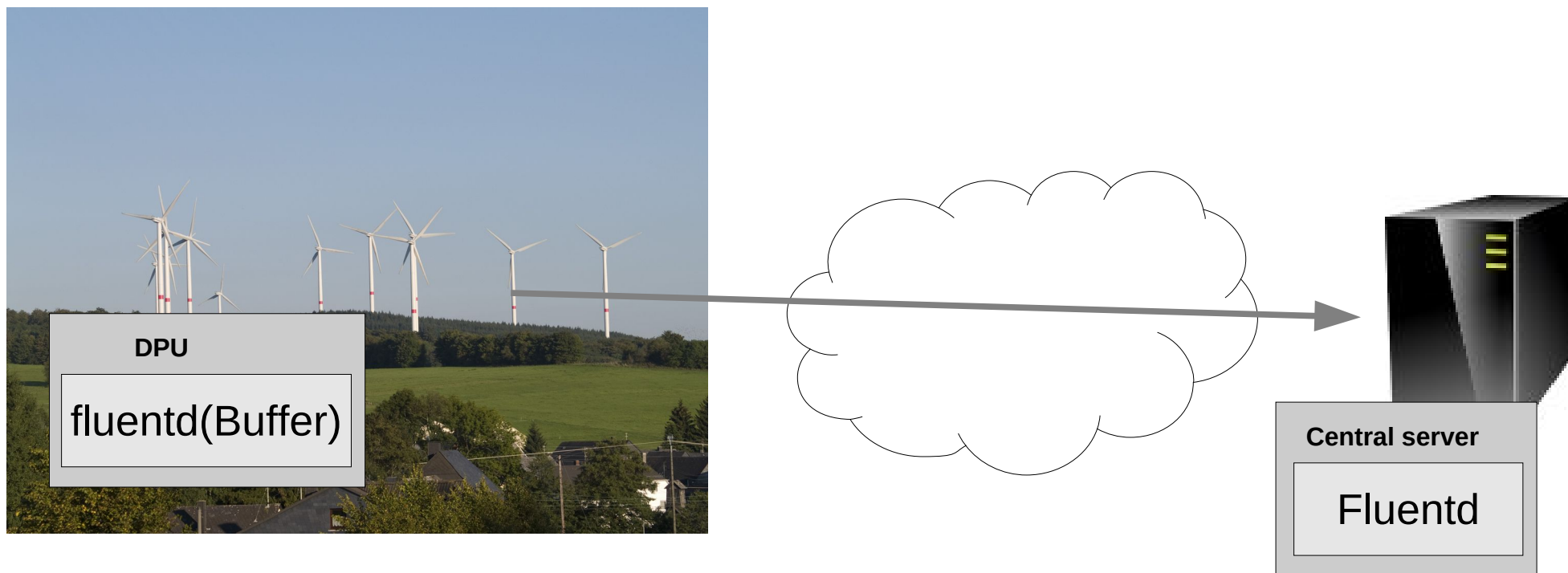
Wendelin Distinctive Advantages

- **Native out-of-core NumPy (scikit-learn, pydata)**
- **Native parallel processing**
- **Bare metal performance (GPU, FORTRAN)**
- **Transactions (ingestion, processing)**
- **NewSQL queries**
- **Built-in PaaS**
- **Lower deployment cost (10x less than...)**

Detailed Example



Data Transportation **fluentd**



3 months benchmark

Frequent downtime (server, network)

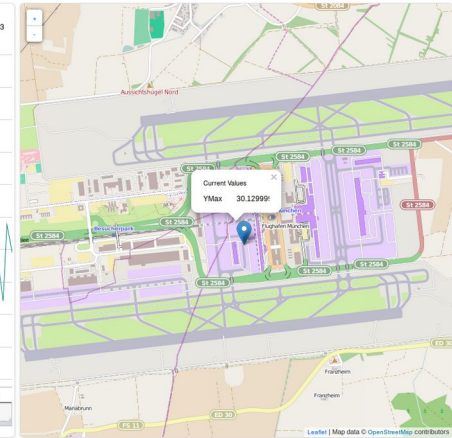
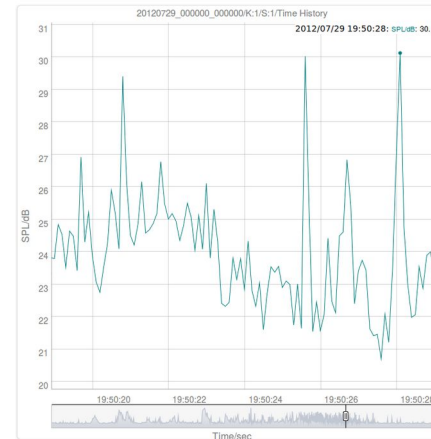
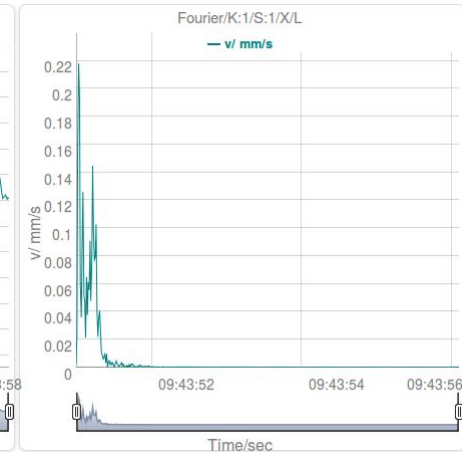
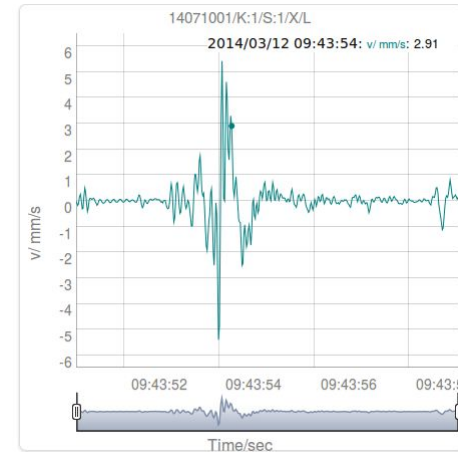
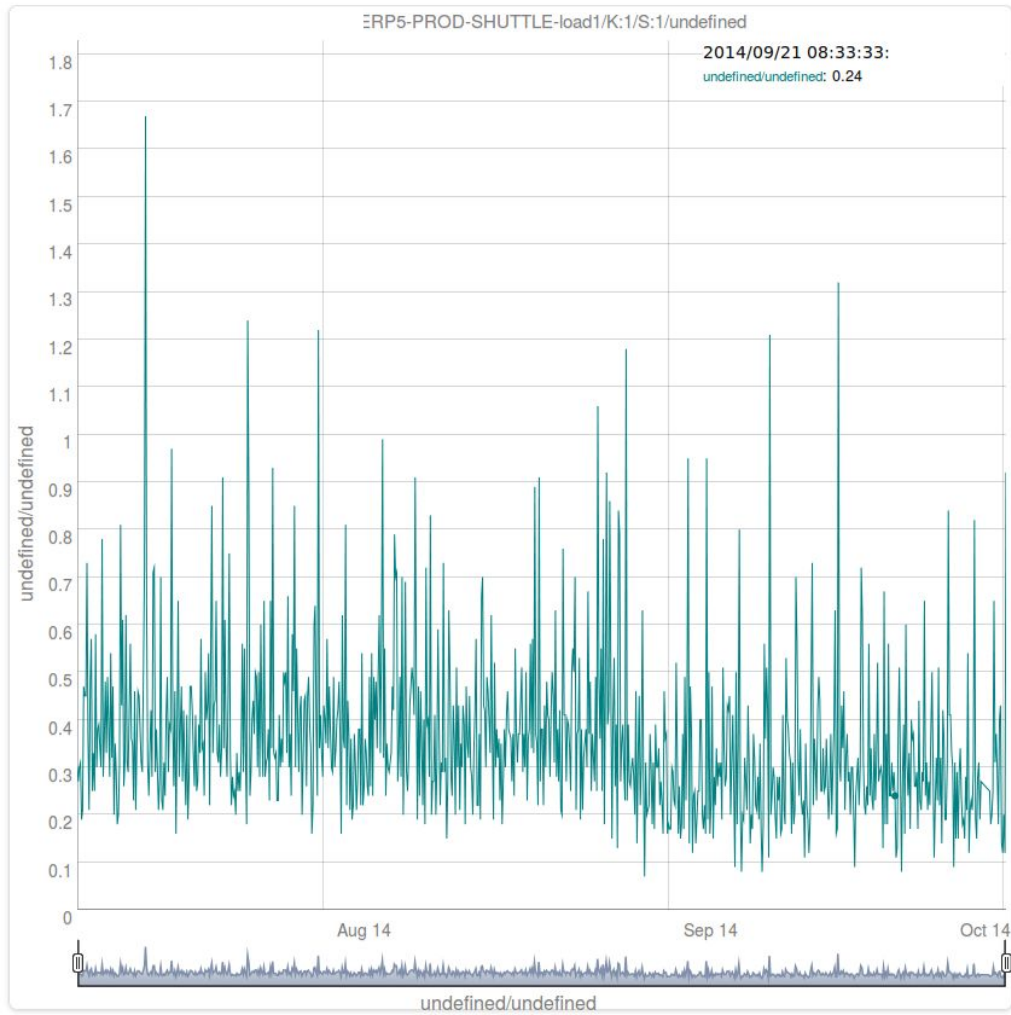
Very poor networking (ADSL, 3G)

< 0.001% loss

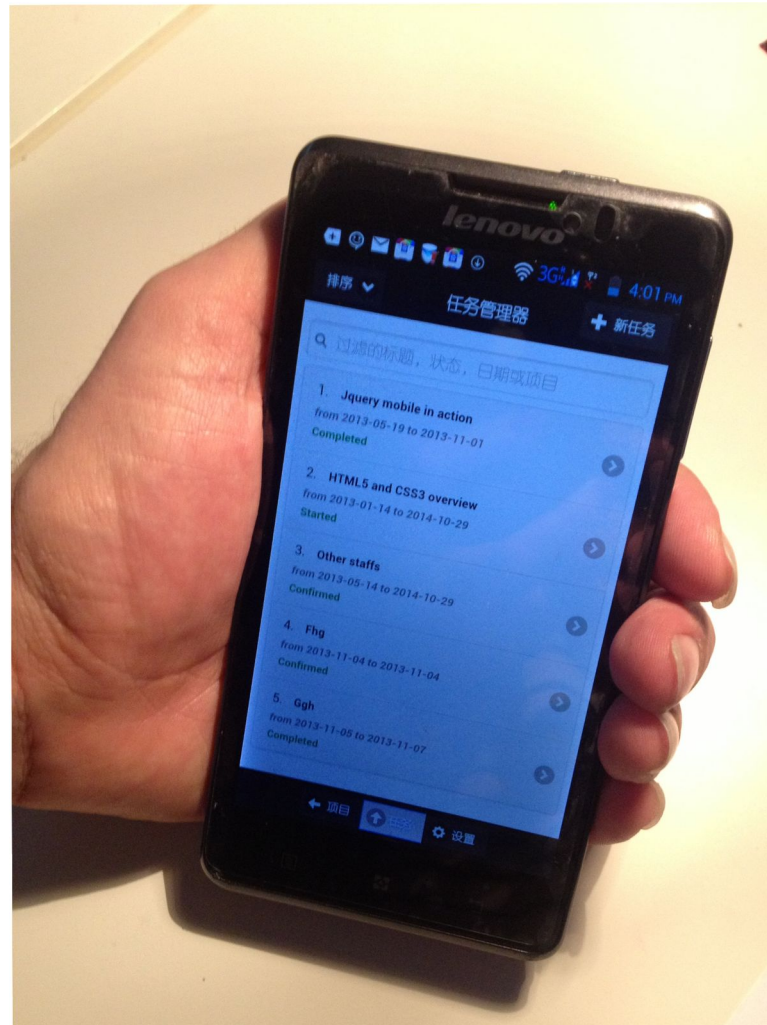
UI: HTML5 Components **RenderJS**



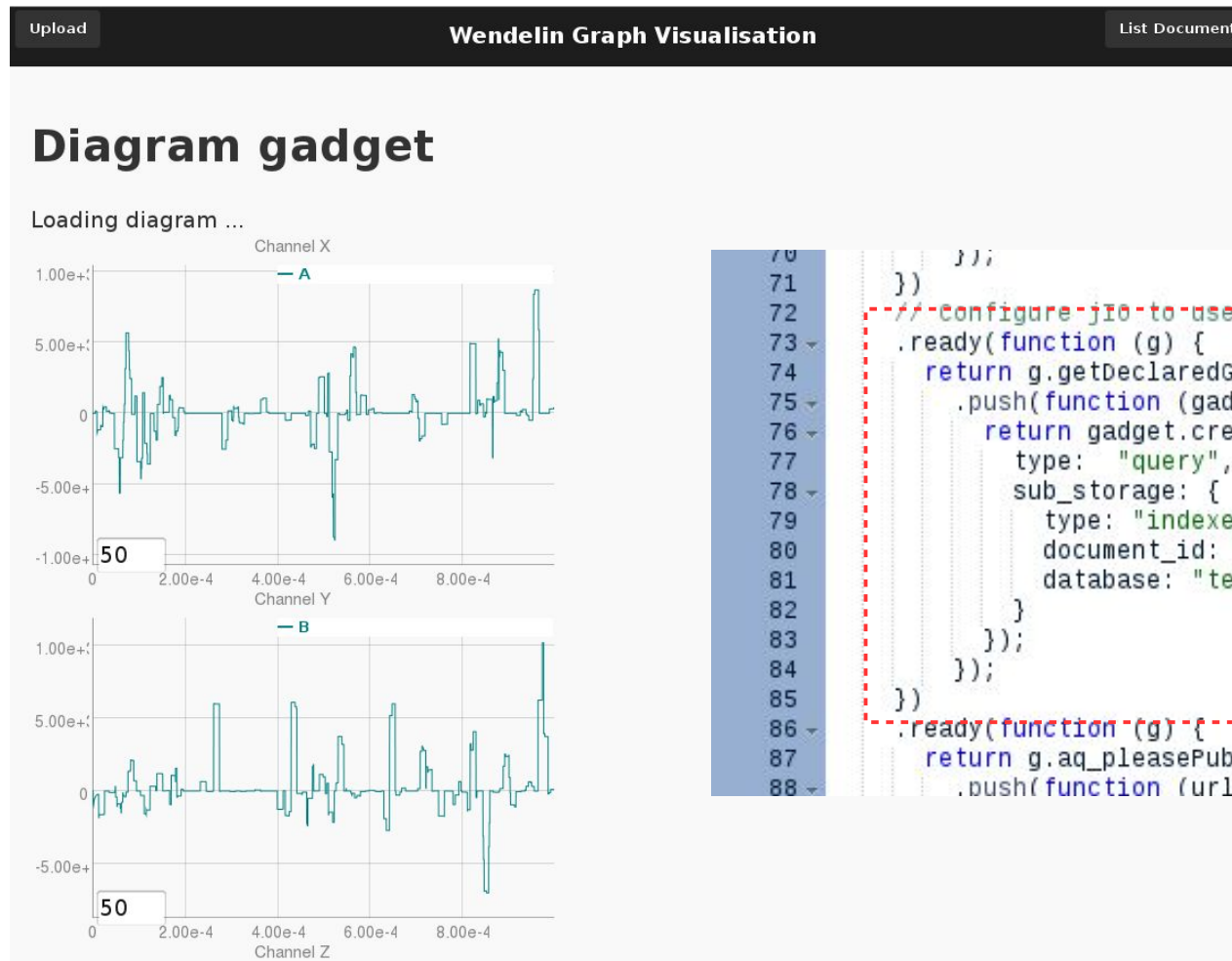
Extend UI Components **RenderJS**



UI : Responsive **RenderJS**

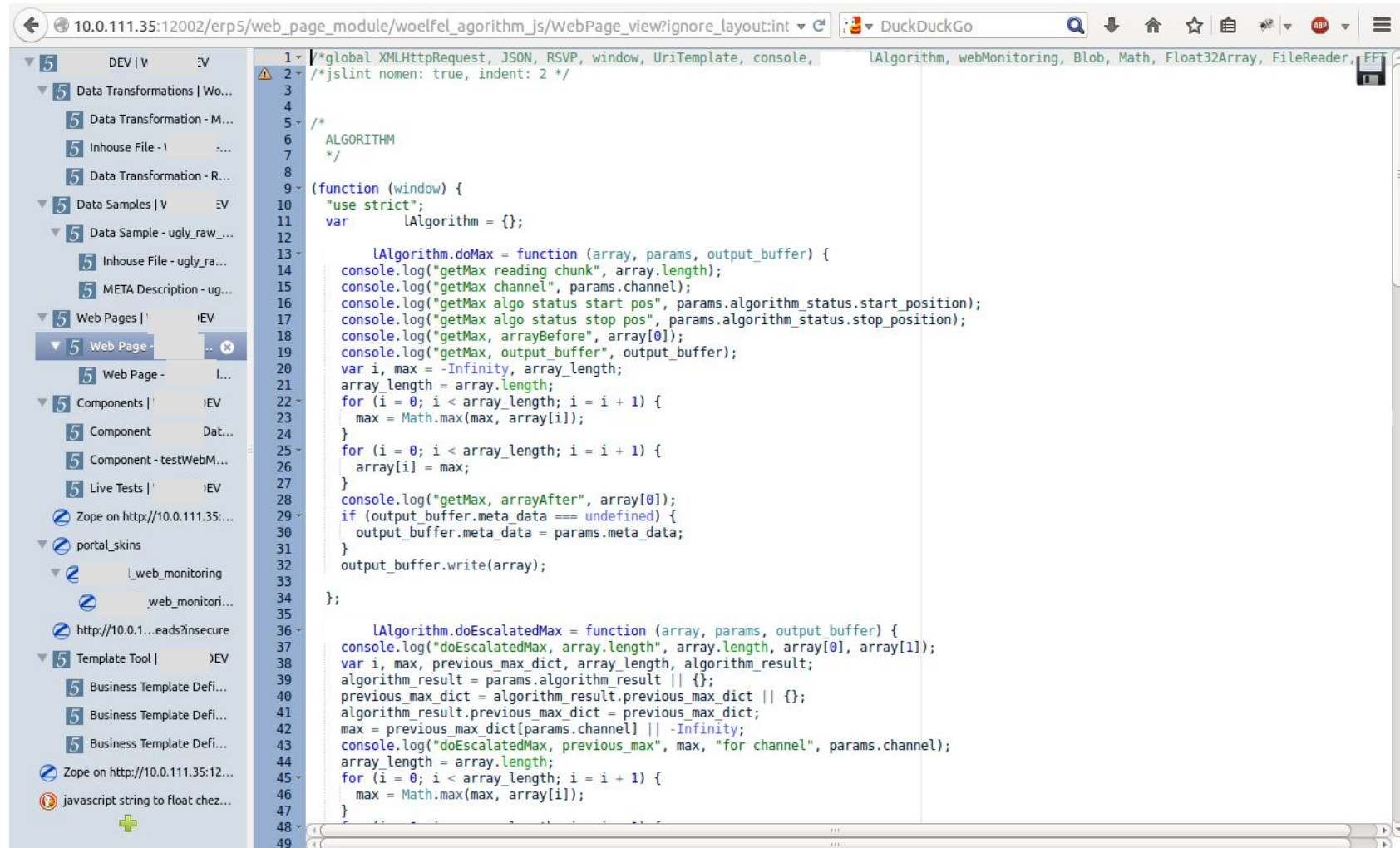


UI: Offline / Other Backends **JIO**



```
70  
71  
72  
73 // Configure jio to use localStorage  
74 .ready(function (g) {  
75   return g.getDeclaredGadget("JIO")  
76   .push(function (gadget) {  
77     return gadget.createJio({  
78       type: "query",  
79       sub_storage: {  
80         type: "indexeddb",  
81         document_id: "/",  
82         database: "test_ivan"  
83       }  
84     });  
85   });  
86 .ready(function (g) {  
87   return g.aq_pleasePublishMyState({page: 'listbox'})  
88   .push(function (url) {
```








Data Science in Javascript vs. Python ?



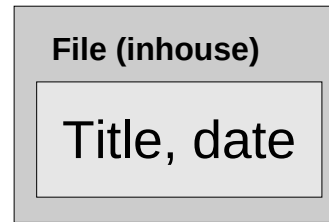
The screenshot shows a web browser window with a URL bar displaying "10.0.111.35:12002/erp5/web_page_module/woelfel_agorithm_js/WebPage_view?ignore_layout:int". The browser's address bar shows "DuckDuckGo". On the left side, there is a file explorer panel showing a directory structure with folders like "Data Transformations", "Data Samples", "Web Pages", and "Components". The main area of the browser displays a JavaScript code editor with the following code:

```
1- /*global XMLHttpRequest, JSON, RSVP, window, UriTemplate, console,
2- /*jslint nomen: true, indent: 2 */
3-
4-
5-
6- /*
7-  ALGORITHM
8- */
9-
10- (function (window) {
11-   "use strict";
12-   var lAlgorithm = {};
13-
14-   lAlgorithm.doMax = function (array, params, output_buffer) {
15-     console.log("getMax reading chunk", array.length);
16-     console.log("getMax channel", params.channel);
17-     console.log("getMax algo status start pos", params.algorithm_status.start_position);
18-     console.log("getMax algo status stop pos", params.algorithm_status.stop_position);
19-     console.log("getMax, arrayBefore", array[0]);
20-     console.log("getMax, output_buffer", output_buffer);
21-     var i, max = -Infinity, array_length;
22-     array_length = array.length;
23-     for (i = 0; i < array_length; i = i + 1) {
24-       max = Math.max(max, array[i]);
25-     }
26-     for (i = 0; i < array_length; i = i + 1) {
27-       array[i] = max;
28-     }
29-     console.log("getMax, arrayAfter", array[0]);
30-     if (output_buffer.meta_data === undefined) {
31-       output_buffer.meta_data = params.meta_data;
32-     }
33-     output_buffer.write(array);
34-   };
35-
36-   lAlgorithm.doEscalatedMax = function (array, params, output_buffer) {
37-     console.log("doEscalatedMax, array.length", array.length, array[0], array[1]);
38-     var i, max, previous_max_dict, array_length, algorithm_result;
39-     algorithm_result = params.algorithm_result || {};
40-     previous_max_dict = algorithm_result.previous_max_dict || {};
41-     algorithm_result.previous_max_dict = previous_max_dict;
42-     max = previous_max_dict[params.channel] || -Infinity;
43-     console.log("doEscalatedMax, previous_max", max, "for channel", params.channel);
44-     array_length = array.length;
45-     for (i = 0; i < array_length; i = i + 1) {
46-       max = Math.max(max, array[i]);
47-     }
48-   };
49- }
```


Data Sciences in Javascript ? **phantomjs**

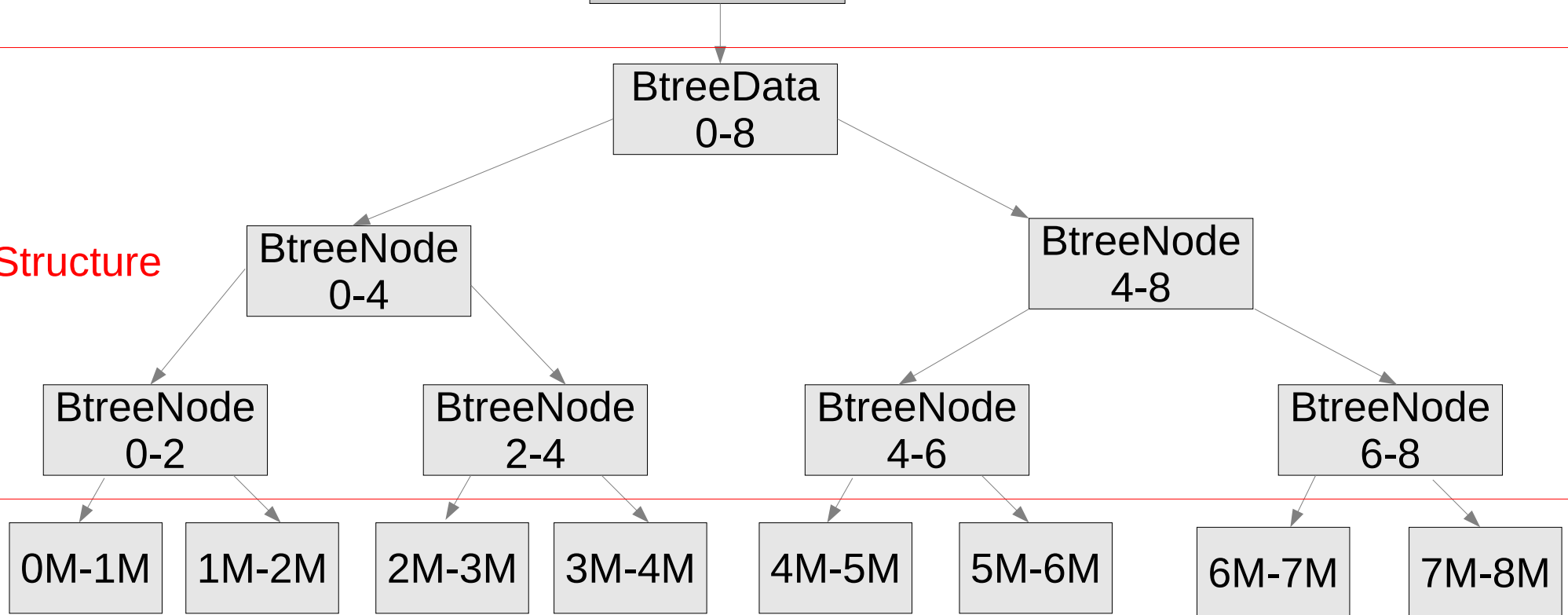
- **Small data on client side** 
- **Small data on server side** 
- **Medium data (> 1 GB) in JS** 
- **Out-of-core data in JS** 
- **PyData compiled in JS** 
- **PyData in NaCl / PNaCl** 

Storing large streams in NEO



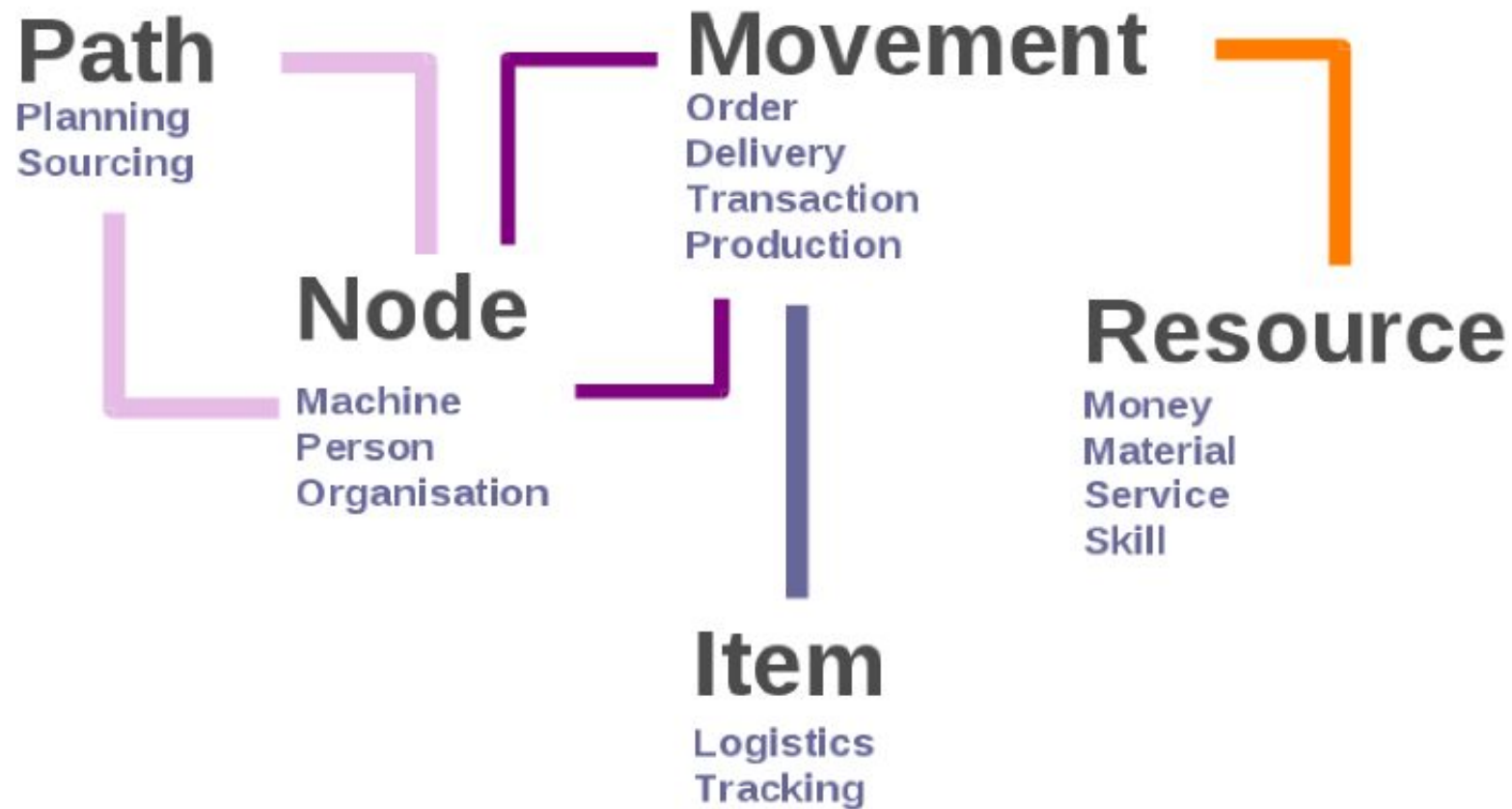
Access: $O(\log(N))$
Overhead : 0.2%

Structure



Data

UBM Monitoring Model?

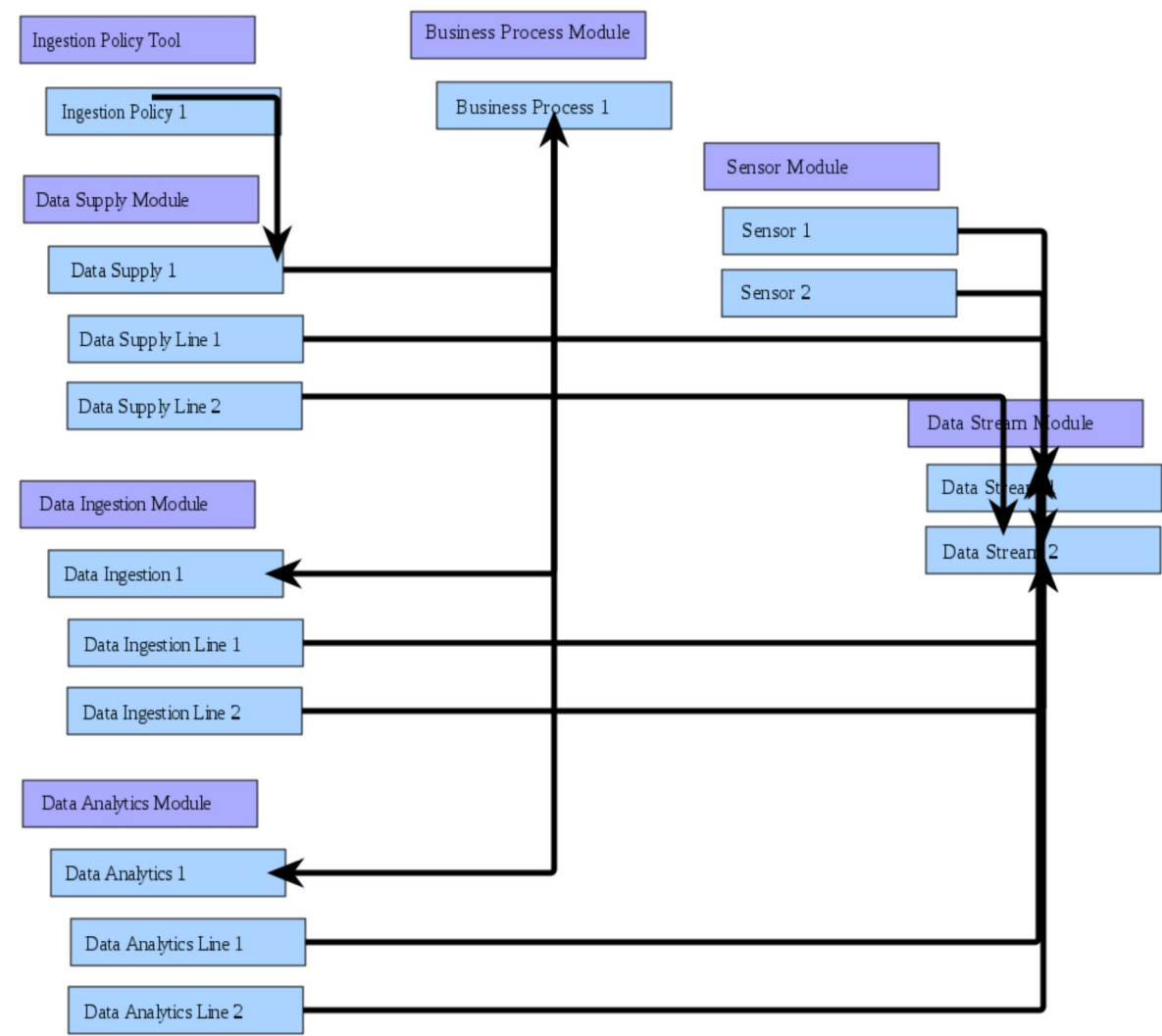


UBM Business Model



- **Movement** – ingestion of data
- **Resource** – type of data (ex. memory log)
- **Node** – data source, data owner
- **Path** – data source registration
- **Item** – sensor, data itself, license, data set

UBM Business Model



What UBM gets us for free

- **Accounting, billing and payment**
- **User registration and management**
- **Rule based security model**
- **Customer relationship management**
- **Web Content Management**

→ **save 12+ months and > 200 K€ on any Big Data project**

Future Roadmap



Roadmap

www.wendelin.io

- **Mainly accelerate learning curve**

- ☐ Universal packaging
- ☐ Ready to use examples
- ☐ Act as a backend to ipython notebook
- ☐ Port joblib to CMFActivity

- **Yet, you can start using part of Wendelin now!**

- ☒ **wendelin.core out-of-core for NumPy**

[PyData Paris 2015 - 16h45 Kirill Smelkov](#)

- ☒ **JIO abstract data access library**

- ☒ **RenderJS components**

<http://learn.renderjs.org>

- ☒ **UI sample application**

<https://lab.nexedi.cn/Tyagov/wendelin/>

- ☒ **Open Source**

R&D Partners

www.wendelin.io

- **Wendelin-IA (FSN)**

- Nexedi
- Abilian
- 2nd Quadrant
- Paris 13
- IMT
- INRIA / ENS
- MMC Rus (Ru)
- X Corp

- **Windelin (Eurostars)**

- Nexedi (FR)
- MariaDB (FI)
- Y Corp (DE)





Wendelin Big Data

Industrial Monitoring Platform

2014-04-03 – Paris



Wendelin Big Data *Industrial Monitoring Platform*

2014-04-03 – Paris

© 2015 Nexedi SA – Company Confidential



Who are we?

- **Jean-Paul Smets**
- **Nexedi CEO**
- **Author of ERP5**
- **jp@nexedi.com**
- **Ivan Tyagov**
- **Senior Developer**
- **Wendelin project lead**
- **ivan@nexedi.com**

Who is missing?

- **Kirill Smelkov**
- **Senior Developer**
- **wendelin.core**
- **Sebastien Robin**
- **Project Director**
- **Author of POC**

Agenda

- **Where do we come from**
- **Wendelin Architecture**
- **Detailed Example**
- **Future Roadmap**

Where do we come from?




© 2015 Nexedi SA – Company Confidential

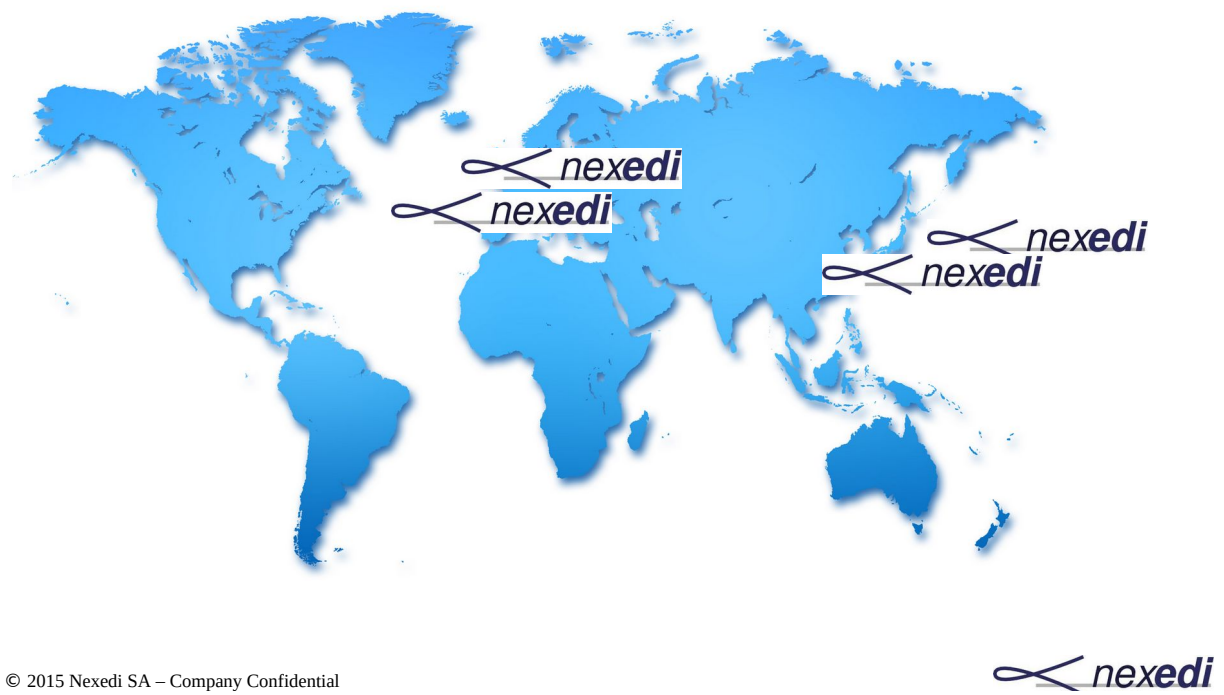


The solution that was deployed at the Lightning Protection Center complies is based on open source software – with full access to source code – and does not use software made by IBM, Oracle or EMC. It is thus a “No IOE” compliant solution, in line with directives published by Chinese governments for certain markets.

Nexedi

- **Possibly Largest OSS Publisher in Europe**

- ERP5: ERP, CRM, ECM, e-business framework
- SlapOS: distributed mesh cloud operation system
- NEO: distributed transactional NoSQL database
-  **Wendelin: out-of-core big data based on NumPy**
- re6st: resilient IPv6 mesh overlay network
- RenderJS: javascript component system
- JIO: javascript virtual database and virtual filesystem
- cloudoon: multimedia conversion server
- Web Runner: web based Platform-as-a-Service (PaaS) and IDE
- OfficeJS: web office suite based on RenderJS and JIO





Aide et Action

L'Education change le monde



WEINPARIS.COM

BY ZHWEE





REPUBLIQUE FRANÇAISE



BCEAO

BAHREIN CENTRAL BANK



CAPAGO

The Schengen link



AIRBUS

DEFENCE & SPACE



Sanef



Kyorin



HANGZHOU DIANZI UNIVERSITY



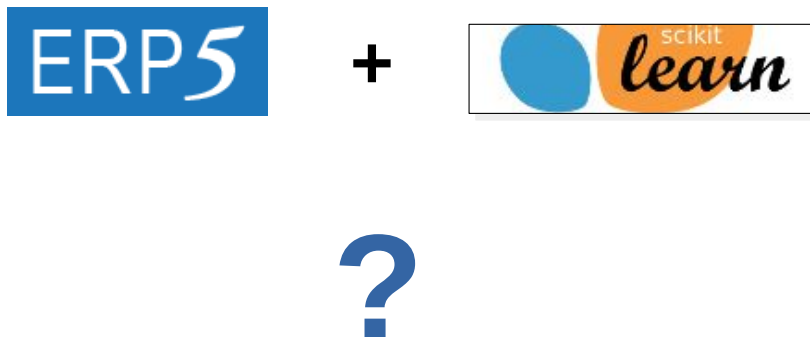
MITSUBISHI
MOTORS



nexedi

© 2015 Nexedi SA – Company Confidential

Application Convergence



Case 1: Wind Turbines



- Collect logs
- Collect records
- Predict failure
- Plan maintenance
- Reduce downtime
 - add X% profits

ERP5The Nexedi logo, which consists of a blue circle with a white dot inside, followed by the word "nexedi" in a stylized, lowercase font.**ERP5**

Case 2: Cars



- Collect logs
- Collect records
- Predict failure
- Plan maintenance
- Reduce downtime
 - increase loyalty

ERP5

learn

ERP5

© 2015 Nexedi SA – Company Confidential

 nexedi

Case 3: Solar Energy



- Collect logs
- Collect records
- Predict degradation
- Plan maintenance
- Increase efficiency
 - add X% to profits



Wendelin Architecture



© 2015 Nexedi SA – Company Confidential



The solution that was deployed at the Lightning Protection Center complies is based on open source software – with full access to source code – and does not use software made by IBM, Oracle or EMC. It is thus a “No IOE” compliant solution, in line with directives published by Chinese governments for certain markets.

Standard Hardware no router / no SAN



inspur 浪潮 | lenovo 联想 | CORETO

x 160

- 2 x 10 Gbps
- 2 x 6 core Xeon CPU
- 512 GB RAM
- 4 x 1 TB SSD
- 1 x M2090 GPU

+



x 32

- 10 Gbps
- Unmanaged

+

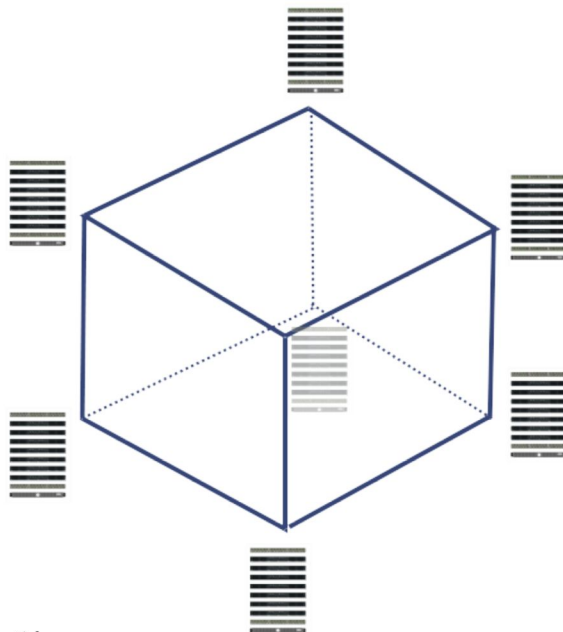


x 320

© 2015 Nexedi SA – Compan

 nexedi

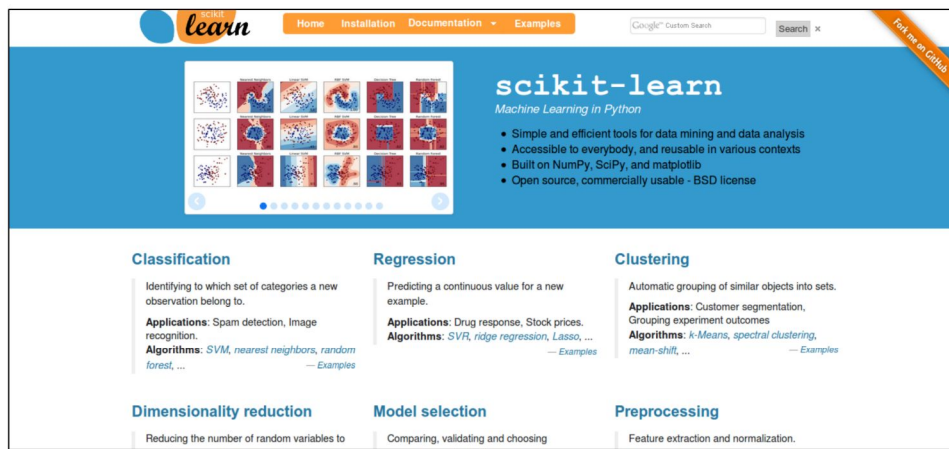
Wendelin Hypercube Datacenter



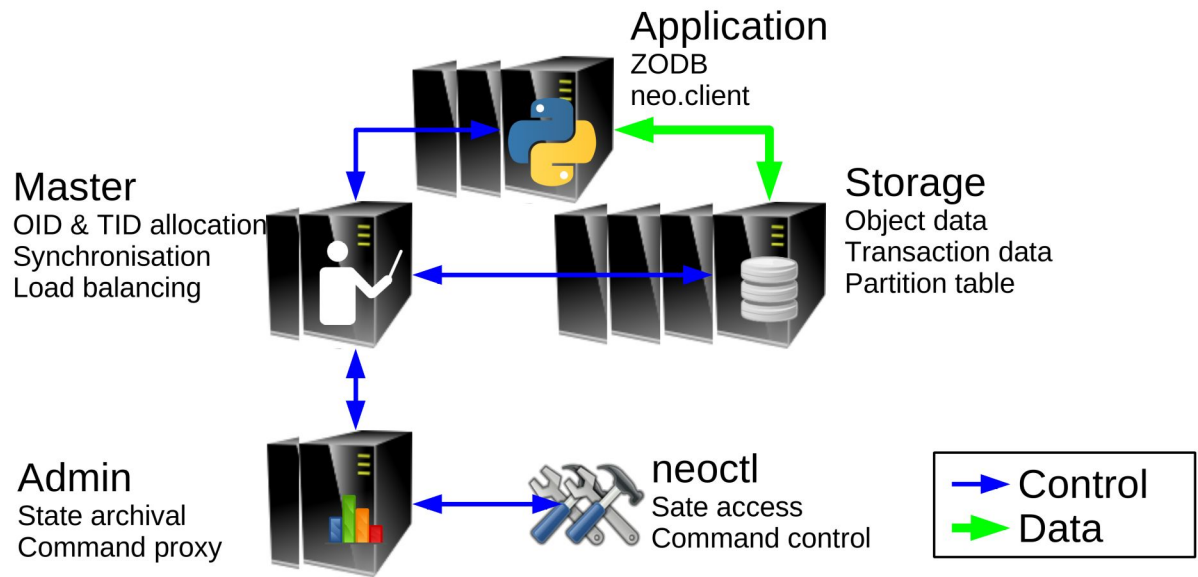
© 2015 Nexedi SA – Company Confidential



Take the Best Analytics scikit-learn.org



Add Distributed Storage neopod.org



“Magic” out-of-core for NumPy

PyData Paris 2015 – 16h45 Kirill Smelkov

ZBigArray

1	2	3	4	5	6	7	8	9	10	11	12
---	---	---	---	---	---	---	---	---	----	----	----



Add Elastic PaaS erp5.com

```
# Initialize data
data_size = 1000000
server_count = 1000
chunk_size = data_size / server_count
data = array(data_size)

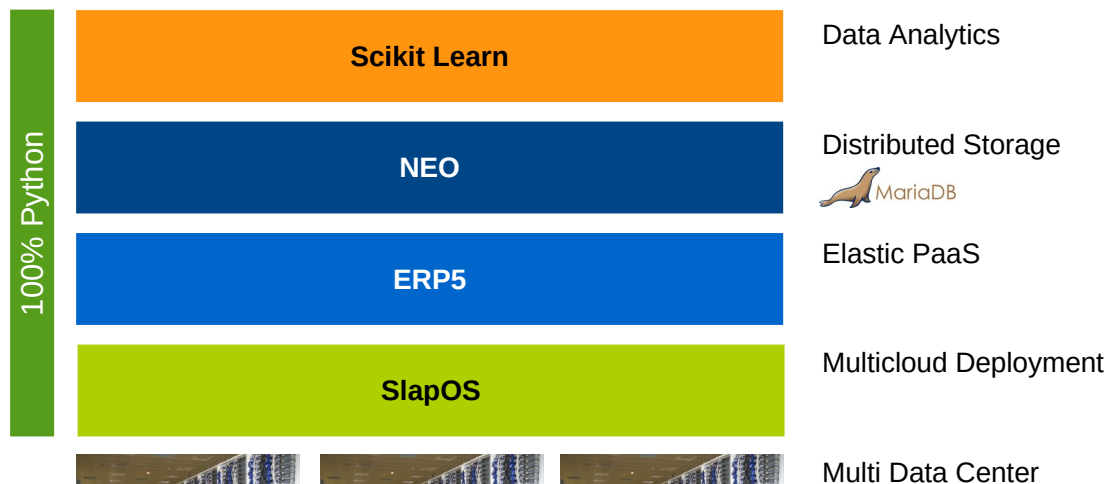
# Process data in parallel on each server (Map Reduce, Batch, etc.)
for server in server_count:
    data.activate().process(server*chunk_size, chunk_size)
```

ERP5 PaaS

And Multicloud Deployment slapos.org



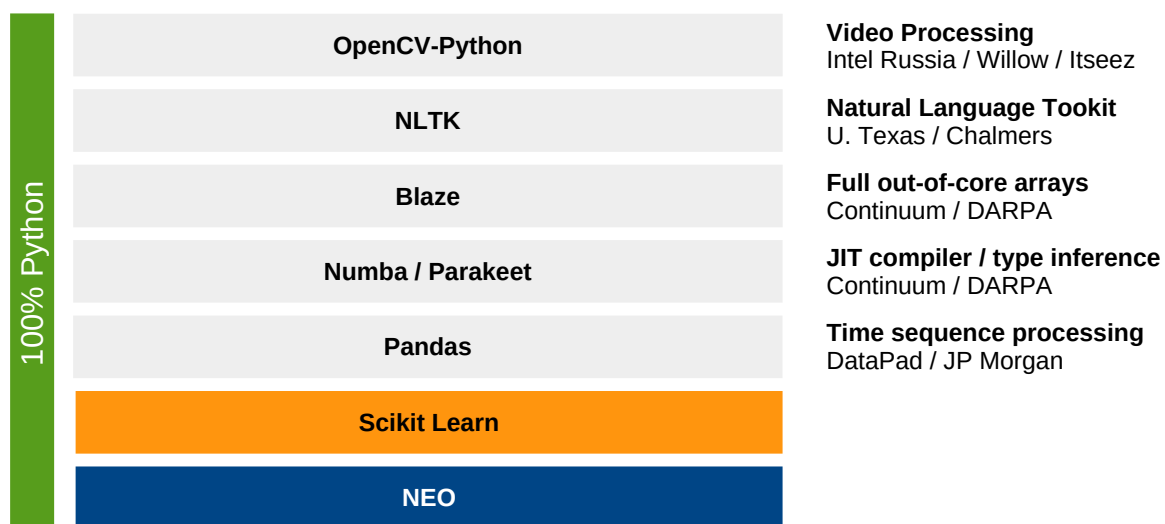
Wendelin Platform 100% open source



© 2015 Nexedi SA – Company Confidential



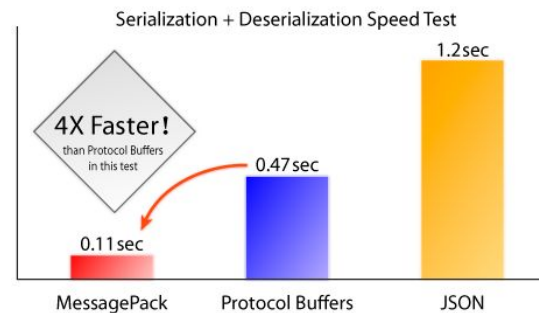
Wendelin Options 100% open source



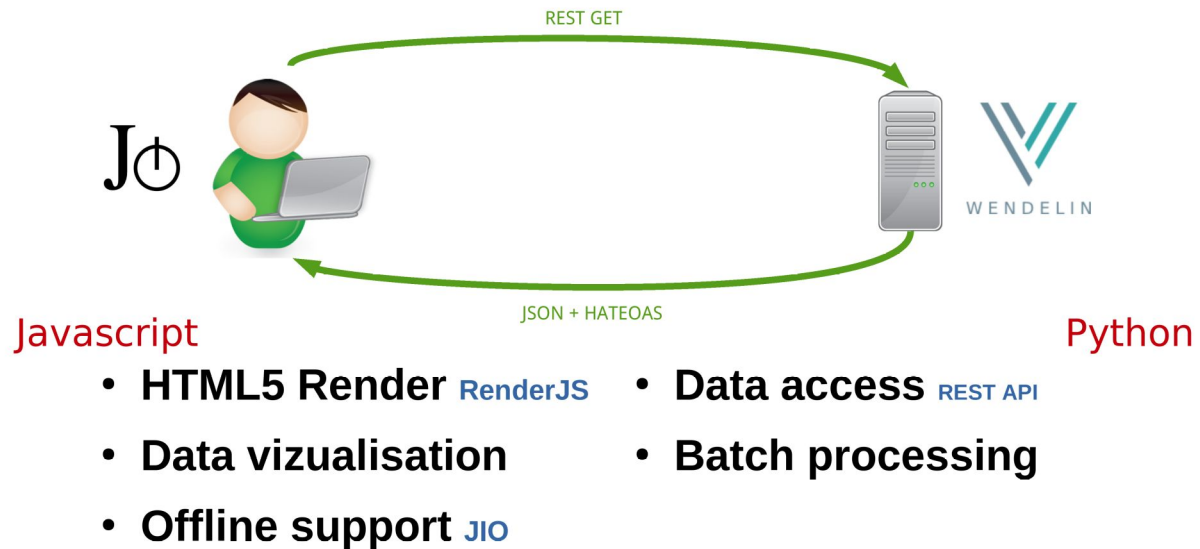
Data Ingestion: fluentd



- **Based on MsgPack middleware**
- **Created by TreasureData (BDaaS pioneers)**
- **Used by Amazon**
- **Numerous plugins**
- **Scalable and resilient**
- **Bandwidth saver**



Wendelin UI



© 2015 Nexedi SA – Company Confidential



The experimental HTML5 UI of ERP5 uses a library called JIO to abstract the relation between the browser and the server.

The browser sends REST requests over HTTP.

The server returns JSON data over HTTP with a self-discoverable format, something called HATEOAS.

The user interface is implemented as a javascript application that runs on the browser side. HTML is generated on the browser side by Javascript code. Form data is prepared by Javascript code and sent as JSON to the server. Many features of the application are still available even offline.

The role of the server in this architecture is only to provide access to the data and to validated the data before updating records in the dabase, using global consistency rules.

Wendelin Distinctive Advantages

- **Native out-of-core NumPy (scikit-learn, pydata)**
- **Native parallel processing**
- **Bare metal performance (GPU, FORTRAN)**
- **Transactions (ingestion, processing)**
- **NewSQL queries**
- **Built-in PaaS**
- **Lower deployment cost (10x less than...)**

Detailed Example

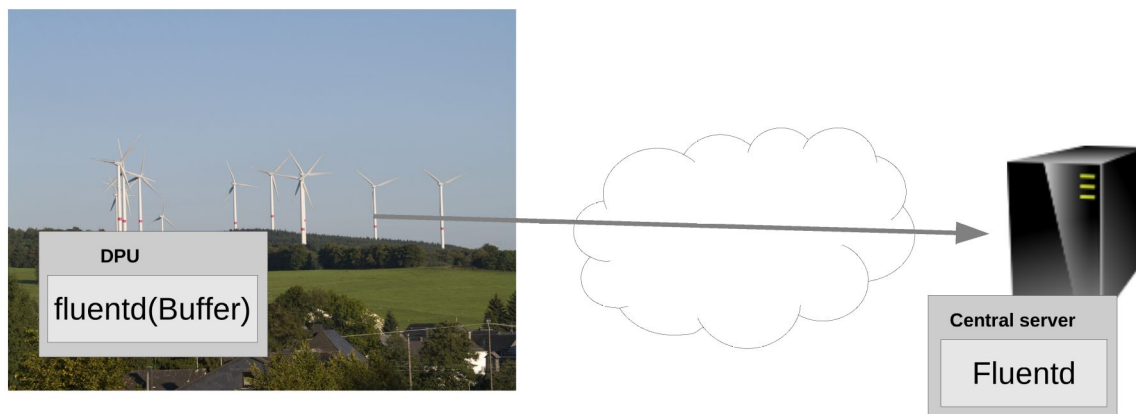


© 2015 Nexedi SA – Company Confidential



The solution that was deployed at the Lightning Protection Center complies is based on open source software – with full access to source code – and does not use software made by IBM, Oracle or EMC. It is thus a “No IOE” compliant solution, in line with directives published by Chinese governments for certain markets.

Data Transportation **fluentd**



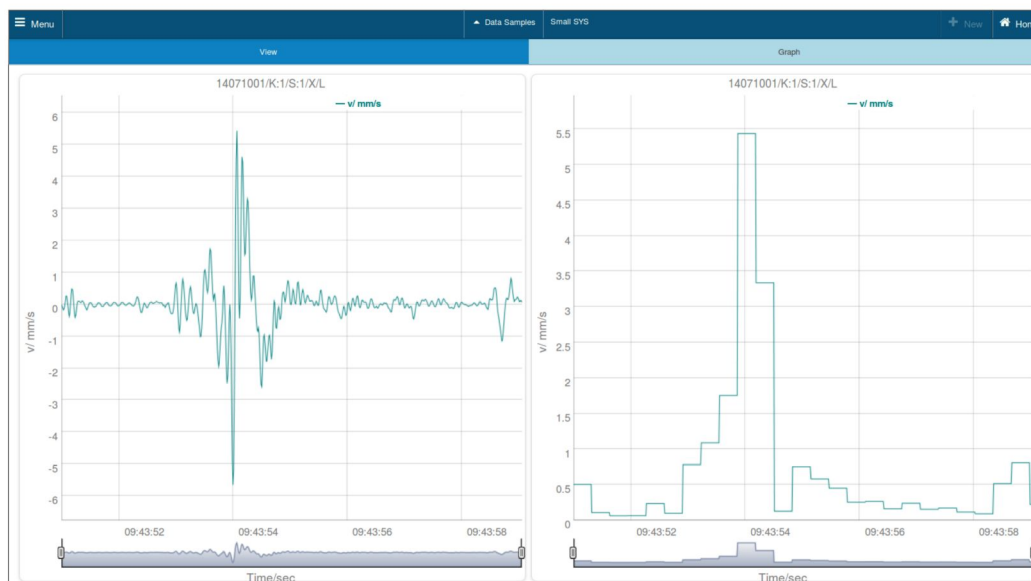
3 months benchmark

Frequent downtime (server, network)

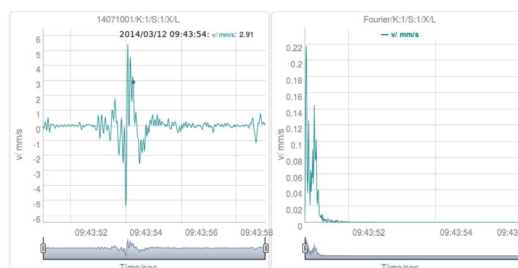
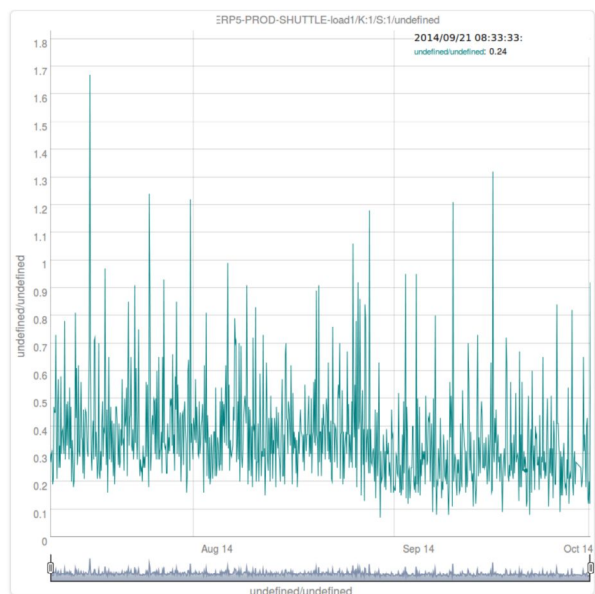
Very poor networking (ADSL, 3G)

< 0.001% loss

UI: HTML5 Components **RenderJS**



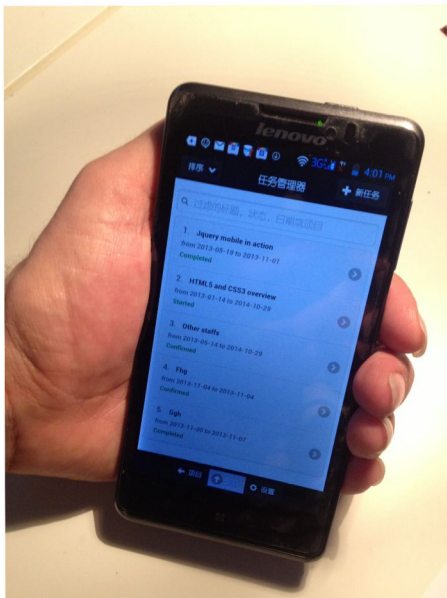
Extend UI Components **RenderJS**



© 2015 Nexedi SA – Company Confidential



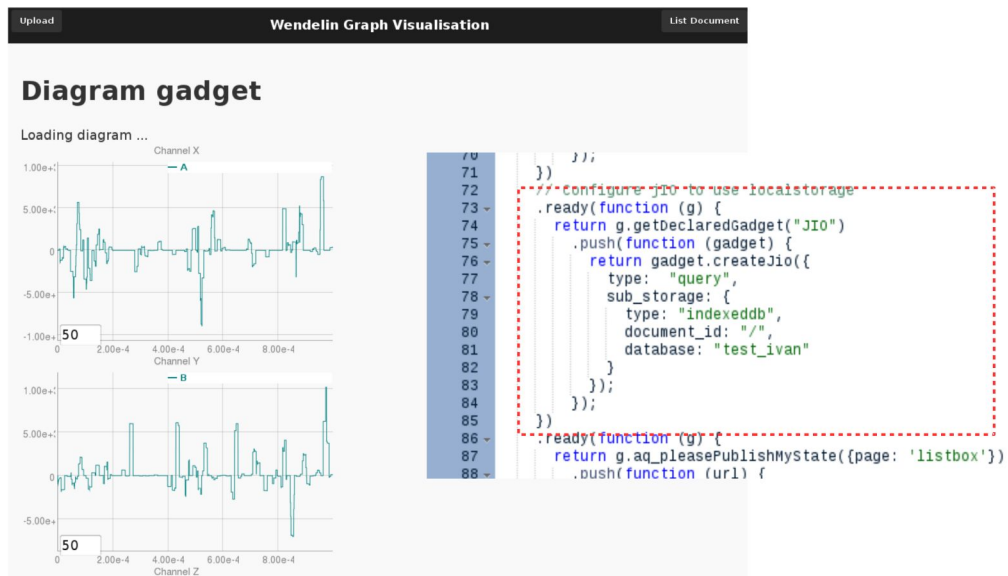
UI : Responsive **RenderJS**



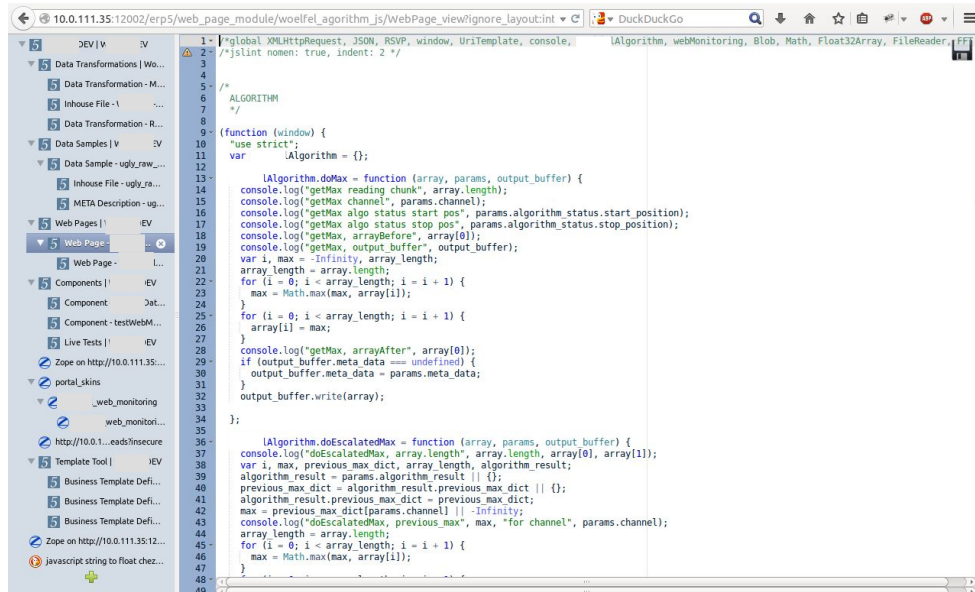
© 2015 Nexedi SA – Company Confidential



UI: Offline / Other Backends **JIO**









Data Science in Javascript vs. Python ?

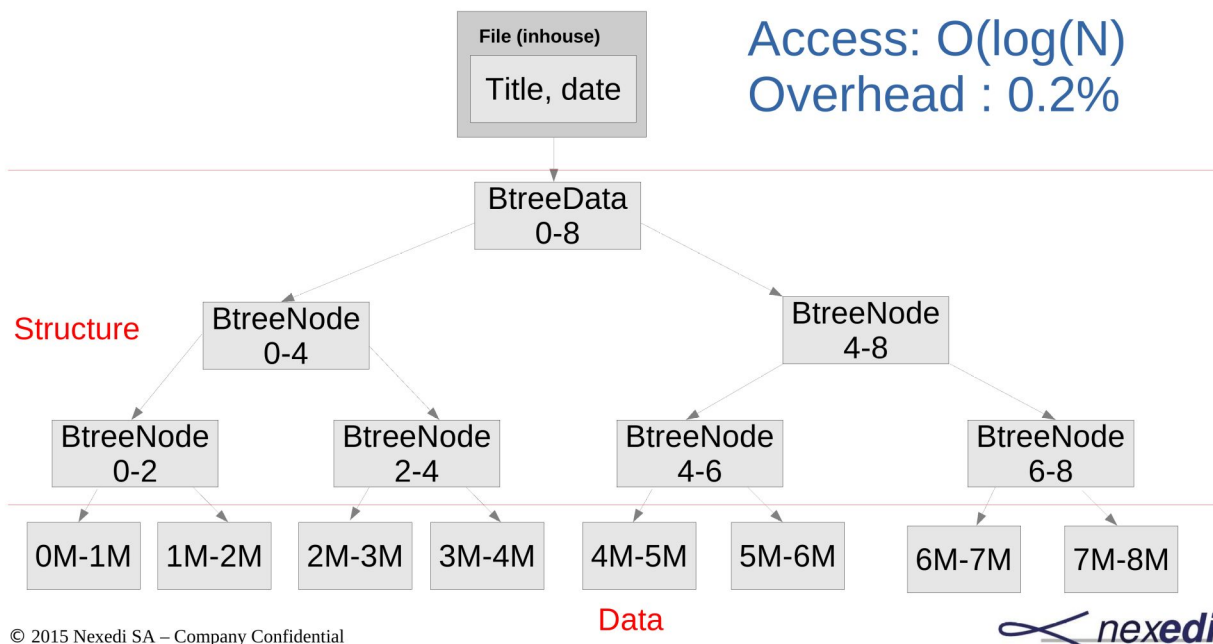


```
1- global XMLHttpRequest, JSON, RSVP, window, UriTemplate, console,
2- /*jslint nomen: true, indent: 2 */
3-
4-
5-
6- /*
7-  * ALGORITHM
8-  */
9-
10- (function (window) {
11-     "use strict";
12-     var
13-         Algorithm = {};
14-
15-     Algorithm.doMax = function (array, params, output_buffer) {
16-         console.log("getMax reading chunk", array.length);
17-         console.log("getMax channel", params.channel);
18-         console.log("getMax algo status start pos", params.algorithm_status.start_position);
19-         console.log("getMax algo status stop pos", params.algorithm_status.stop_position);
20-         console.log("getMax arrayBefore", array[0]);
21-         console.log("getMax output_buffer", output_buffer);
22-         var i, max = -Infinity, array_length;
23-         array_length = array.length;
24-         for (i = 0; i < array_length; i = i + 1) {
25-             max = Math.max(max, array[i]);
26-         }
27-         for (i = 0; i < array_length; i = i + 1) {
28-             array[i] = max;
29-         }
30-         console.log("getMax arrayAfter", array[0]);
31-         if (output_buffer.meta_data === undefined) {
32-             output_buffer.meta_data = params.meta_data;
33-         }
34-         output_buffer.write(array);
35-     };
36-
37-     Algorithm.doEscalatedMax = function (array, params, output_buffer) {
38-         console.log("doEscalatedMax array.length", array.length, array[0], array[1]);
39-         var i, max, previous_max_dict, array_length, algorithm_result;
40-         algorithm_result = params.algorithm_result || {};
41-         previous_max_dict = algorithm_result.previous_max_dict || {};
42-         algorithm_result.previous_max_dict = previous_max_dict;
43-         max = previous_max_dict[params.channel] || -Infinity;
44-         console.log("doEscalatedMax previous_max", max, "for channel", params.channel);
45-         array_length = array.length;
46-         for (i = 0; i < array_length; i = i + 1) {
47-             max = Math.max(max, array[i]);
48-         }
49-     };
50- }
```

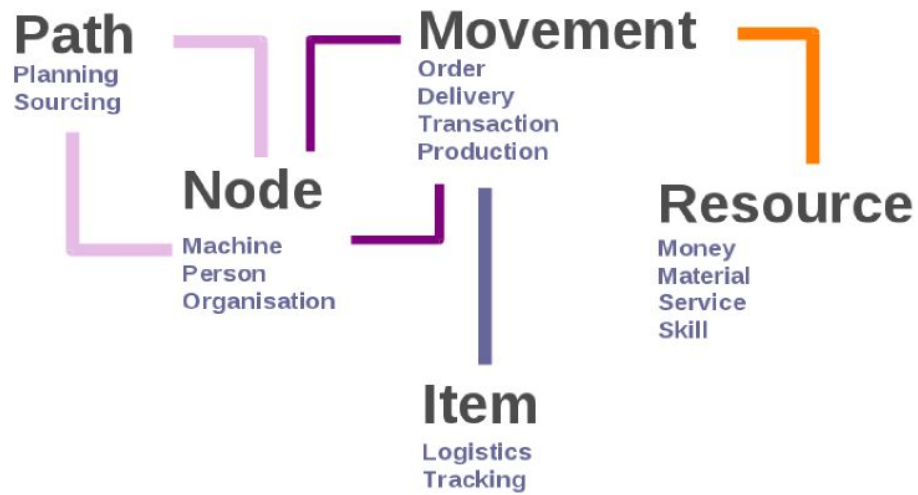
Data Sciences in Javascript ? [phantomjs](#)

- Small data on client side 
- Small data on server side 
- Medium data (> 1 GB) in JS 
- Out-of-core data in JS 
- PyData compiled in JS 
- PyData in NaCl / PNaCl 

Storing large streams in NEO



UBM Monitoring Model?

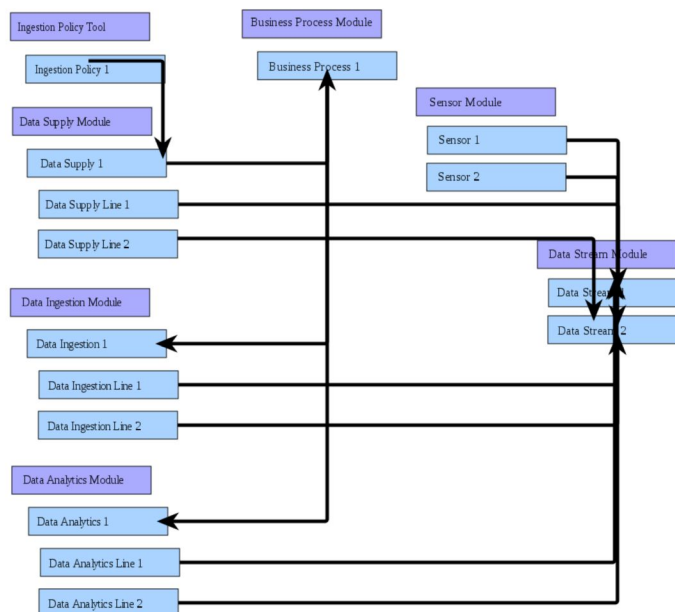


UBM Business Model



- **Movement** – ingestion of data
- **Resource** – type of data (ex. memory log)
- **Node** – data source, data owner
- **Path** – data source registration
- **Item** – sensor, data itself, license, data set

UBM Business Model



What UBM gets us for free

- Accounting, billing and payment
- User registration and management
- Rule based security model
- Customer relationship management
- Web Content Management

→ **save 12+ months and > 200 K€ on any Big Data project**

Future Roadmap



© 2015 Nexedi SA – Company Confidential



The solution that was deployed at the Lightning Protection Center complies is based on open source software – with full access to source code – and does not use software made by IBM, Oracle or EMC. It is thus a “No IOE” compliant solution, in line with directives published by Chinese governments for certain markets.

Roadmap

www.wendelin.io

- **Mainly accelerate learning curve**

- ☐ Universal packaging
- ☐ Ready to use examples
- ☐ Act as a backend to ipython notebook
- ☐ Port joblib to CMFActivity

- **Yet, you can start using part of Wendelin now!**

- ☒ **wendelin.core out-of-core for NumPy** [PyData Paris 2015 - 16h45 Kirill Smelkov](#)
- ☒ **JIO abstract data access library**
- ☒ **RenderJS components** <http://learn.renderjs.org>
- ☒ **UI sample application** <https://lab.nexedi.cn/Tyagov/wendelin/>
- ☒ **Open Source**

R&D Partners

www.wendelin.io

- **Wendelin-IA (FSN)**

- Nexedi
- Abilian
- 2nd Quadrant
- Paris 13
- IMT
- INRIA / ENS
- MMC Rus (Ru)
- X Corp



- **Windelin (Eurostars)**

- Nexedi (FR)
- MariaDB (FI)
- Y Corp (DE)



Wendelin Big Data *Industrial Monitoring Platform*

2014-04-03 – Paris

© 2015 Nexedi SA – Company Confidential

