

Quickstart Big Data

Dr. Olaf Flebbe
of ät oflebbe.de

Fosdem 4.2.2017

HPC, Big Data and Data Science devroom

About me

PhD in computational physics

Former projects: Minix68k (68k FP Emulation), Linux libm.so.5 (High Precision FP), perl and python for epoc, flightgear, msktutil...

Member of PMC of Apache Bigtop

Software Architect for connected E-Bikes

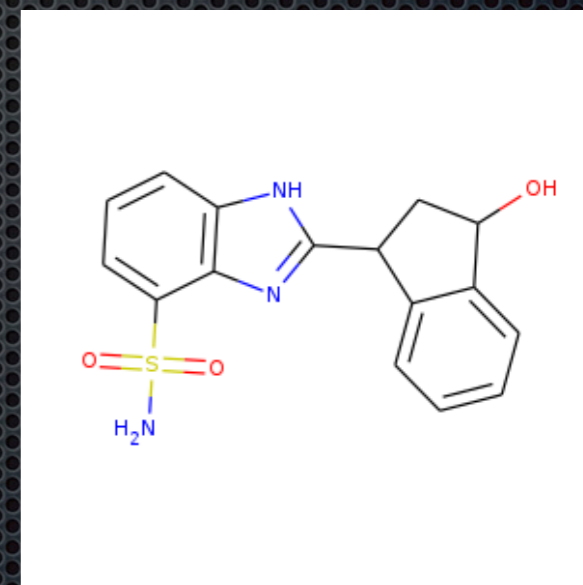


Problem from ChemInformatics

- ✦ Database
For instance:
Synthetically Accessible Virtual Inventory (SAVI) project

- ✦ SMILES Notation

```
OC1CC(C2=CC=CC=C12)C4=NC3=C(C=C  
C=C3[S](=O)(=O)N)[NH]4
```



Problem: Substructure Search

- Commercial solution: Enterprise DB and a Cartridge
- Look out for a Big Data Approach
- Use python rdkit

```
from rdkit import Chem
```

```
m = Chem.MolFromSmiles('OC1CC(C2=CC=CC=C12)C4=NC3=' +  
'C(C=CC=C3[S](=O)(=O)N)[NH]4')
```

```
m.HasSubstructMatch(Chem.MolFromSmiles('C[C@H](F)Cl'))
```

Problem at Hand

- ✦ Ingredients
 - ✦ Time consuming python code
 - ✦ Large Datasets: Even more time consuming
- ✦ Environment
 - ✦ Big Data Cluster
 - ✦ HPC Cluster

How (not) to scale out

- ✦ Most of the runtime is needed for constructing the molecule object
 - ✦ AFAIK One of the problem is to identify ring structures (to aromatize)
- ✦ Can be improved by pickling (serializing) Molecule Objects ...
 - ✦ A huge gain but not scaling effect

How to scale out

- Problem at hand is EP (embarrassingly parallel)
- Can be solved by "Distribute the Algorithm to the data"
- Use Apache Spark-Core python binding
 - Use the RDD paradigm
 - Runs in HPC and BigData Environments

How to scale out

- Read the Instructions at Spark.a.o

```
def matchmol( smile) :  
    return  
Chem.MolFromSmiles( smile).\br/>HasSubstructMatch('C[C@H](F)Cl')  
  
input = sc.textFile("file:///directory/  
smiles.txt")  
input.filter( matchmol).count()
```


HPC Mode

- ✦ HPC Cluster of Machines with
 - ✦ Use Cluster Filesystem for
 - ✦ Deploying Apache Spark
 - ✦ Distribute Data
 - ✦ Does not use locality of data, but works.
- ✦ Use for instance Standalone Mode

Big Data Setup

- Well

Big Data Distro

```

:ass..
=X522nai,>_
=n-- +!!!"^^
.vX>
.)e<o;.
.v2`-{S>
..<de~..;)Sa,
.._aoX}:===>=-?Xo>,
.._aaoZe!`=><i=s+s;~*XXos,,
.....=iisaaoXXZY!"~._v(d=:nc-1s,~?SX#Xouass,,.....
=XXoXXXSXXXXXXXXXZUX21?!"^^_au*`=u2` ]X>.+*a>,-"!Y1XSSX##ZZXXXXXXXXXXXXoXXc
.{XXXXXX2*?!!"^^----_aa2!^-=dX(.+XXc.~!1nas,,----~^"!"!!!?YSXXXXX2+
-"YSXXxo=. _=ssaaav1!!~_aXxe` )SXo>. ~"?Yoouass_s,, _vXXXX2}~
-{XXZoai%%*XXSSSX>.. <uXXX2~ {XXXXs,, =dXXXZX2lii%uXXXXe-
.<XXXXX%- <XXXXX1|==%vdXXXXXo;:.. _vXXXXXXos_ =i|*XXXXX> -<XXXXX`
=5XXXZc .nXXX2> ---=2XXX2^-"|||}"--~{ZXXX1-- :XXXXo; .)XXXX2`
=XXXXZc .nXXX> =XXXe._s=>...)XXXX1 .:SXXXo; .)XXXX2..
<XXXXXc .nXXXS> =XXXosummmmBmma,)ZXXX1 :XXXX2; .)XXXXX.
.<XXXXX( :nXXXS; <XXXXXm#mmmWmmmmmoZXXX1 .3XXXo; .)XXXXX;
.nXXXXX; :XXXXX; =XXXXXmmmBmmWmB#XXXXX1 .nXXXX> :XXXXXc
=oZXXXe; <XXXX2` .)XXXXZmBmBmWmmW#2XXXX1 .vXXXXc vXXXXo;
+Y3S2Xz__..vXXXXe .)ZXXXZmmWmBmBBm#XXXXXo.. {XXXXz:.._vSS2Y1=
---+"*!*Y1s|===uSSSXZUXUXUXUXS2XX2n|_ =| |%Y*?!"^^----
-----
.o. 0000 0000000000. 080
.888. `888 `888' `Y8b `"' .o8
.8"888. 00.00000. .0000. .00000. 888 .00. .00000. 888 8880000 .00000000.o88800 .00000. 00.00000.
.8' `888. 888' `88b`P )88b d88' `Y8 888P`Y88b d88' `88b 8880000888'`888 888' `88b 888 d88' `88b 888' `88b
.880008888. 888 888 .oP"888 888 888 888 888000888 888 `88b 888 888 888 888 888 888 888 888
.8' `888. 888 888d8( 888 888 .o8 888 888 888 .o 888 .88P 888 `88bod8P' 888 .888 888 888 888
o88o o8888o 888bod8P'`Y888"8o`Y8bod8P'o888o o888o`Y8bod8P' o888bood8P' o888o`8000000. "888"Y8bod8P' 888bod8P'
888 d" YD 888
o888o "Y88888P' o888o

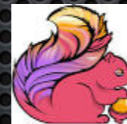
```

Apache Bigtop

- Apache Bigtop is the „Debian“ of the Big Data Distributions
 - reused by Google for their Managed Hadoop Service
 - reused within Cloudera and Hortonworks
 - used by Canonicals Hadoop Offering
 - reused by the ODPI.org



Some components of Apache Bigtop



Components

- Compile Environment (based on docker)
- Convenience artifacts (i.e. repositories for Centos7, Centos6, Debian 8, Ubuntu 16.04, Ubuntu 14.04, Fedora 20, opensuse 42.1)
- Provisioning with docker compose (openstack)
- Deployment Templates based on puppet
- Orchestration with Juju Charms
- Automatic Testing Environment: iTests
- And ... non intel architectures (ppc64le, aarch64)

Use with Bigtop

- Use puppet scripts to deploy Apache Hadoop
 - Use HDFS
 - Use YARN Mode
- Download and unpack Spark 2.1.x

Bigtop 1.2 (preview)



- Will have Spark 2.1
- OpenJDK 8 Support
- Not finished ;-(
- We need help ...
- join at bigtop.apache.org

Conclusion

- Problem runs much better in HPC environment
 - because it is compute intensive
 - not limited by data pipeline bandwidth
- Scales really well $O(n)$
- HW amount needed to challenge the existing solution is too large... need to investigate further

Thanks

of ät oflebbe dot de