

O'REILLY®


Strata


CONFERENCE

—+—

HADOOP

 **WORLD**

 Oct. 23–25, 2012

 NEW YORK, NY

Co-presented by

O'REILLY® **cloudera**

Best Practices for Reproducible Research: Vignettes in Quant Finance

Chang She
@changhiskhan

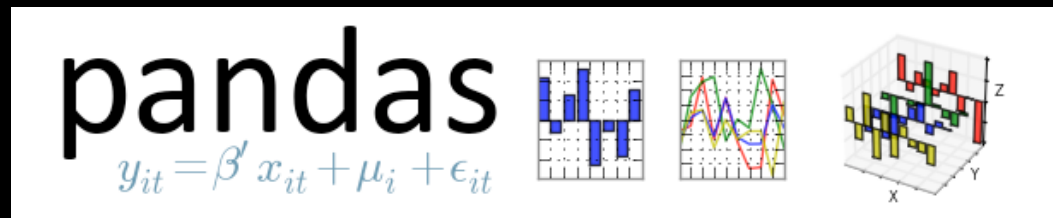
 **LAMBDA·FOUNDRY**

Outline

- Why reproducibility matters
- What is quantitative finance
- Reproducible research → Organize
Version → Code
Test Configurations
Data
- Conclusions

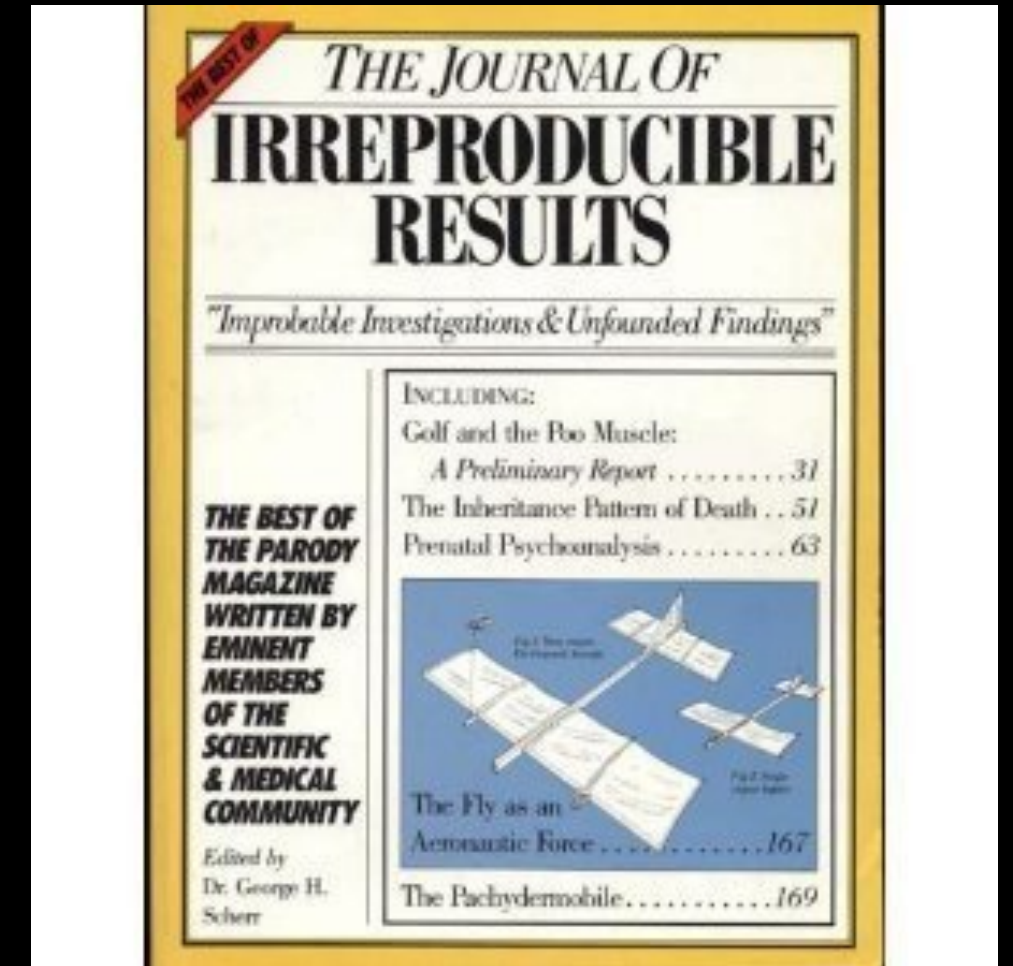
About me

- MIT '05, quant equities/FX 2006-2012
- LF, data tooling, financial analytics

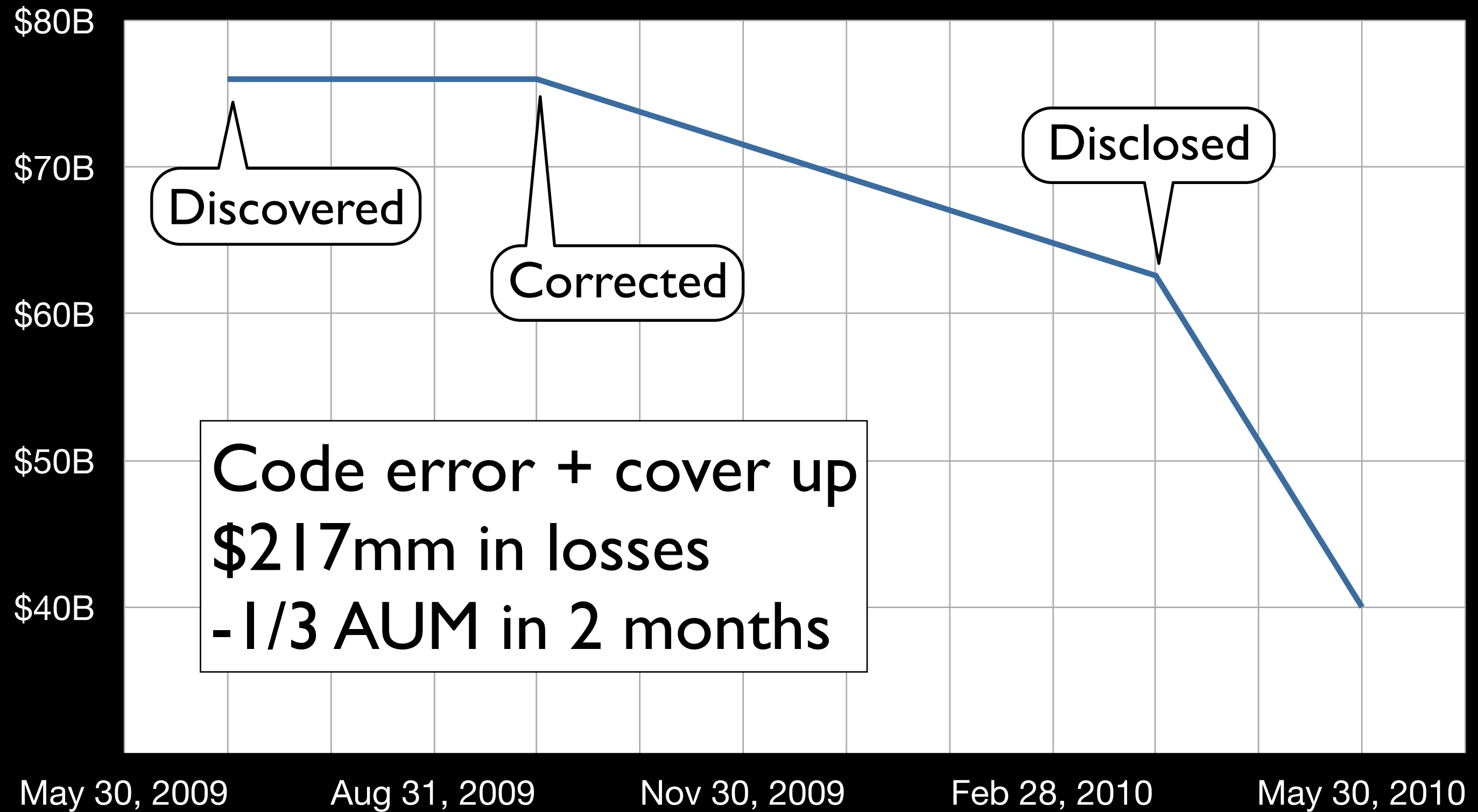


Data Science

- Reproducibility is a cornerstone of modern science
- Not just an academic exercise
- Big data => reproducibility is an organizational effort



AXA Rosenberg AUM



O'REILLY

Strata
CONFERENCE

+ **HADOOP**
WORLD

Relevant questions

- Are my results correct?
- Can anyone produce the same results?
- How do results differ across environments?
- How quickly can I attribute breakages?

Quantitative finance (by marketing)

2: Magic Sauce

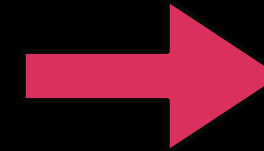
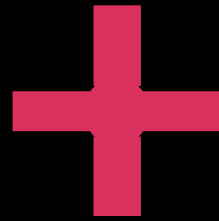
1: Get Data

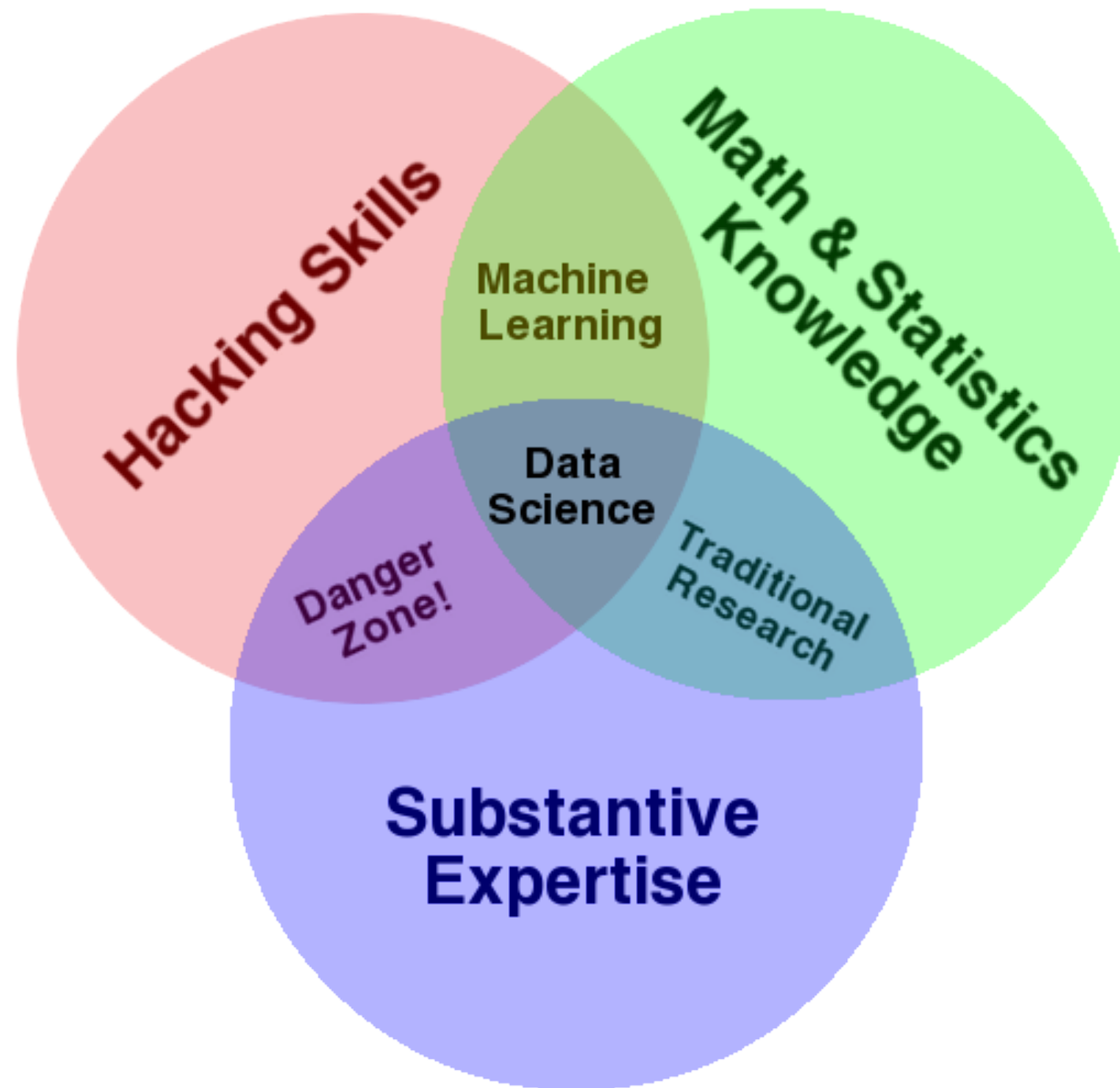


3: Profit!

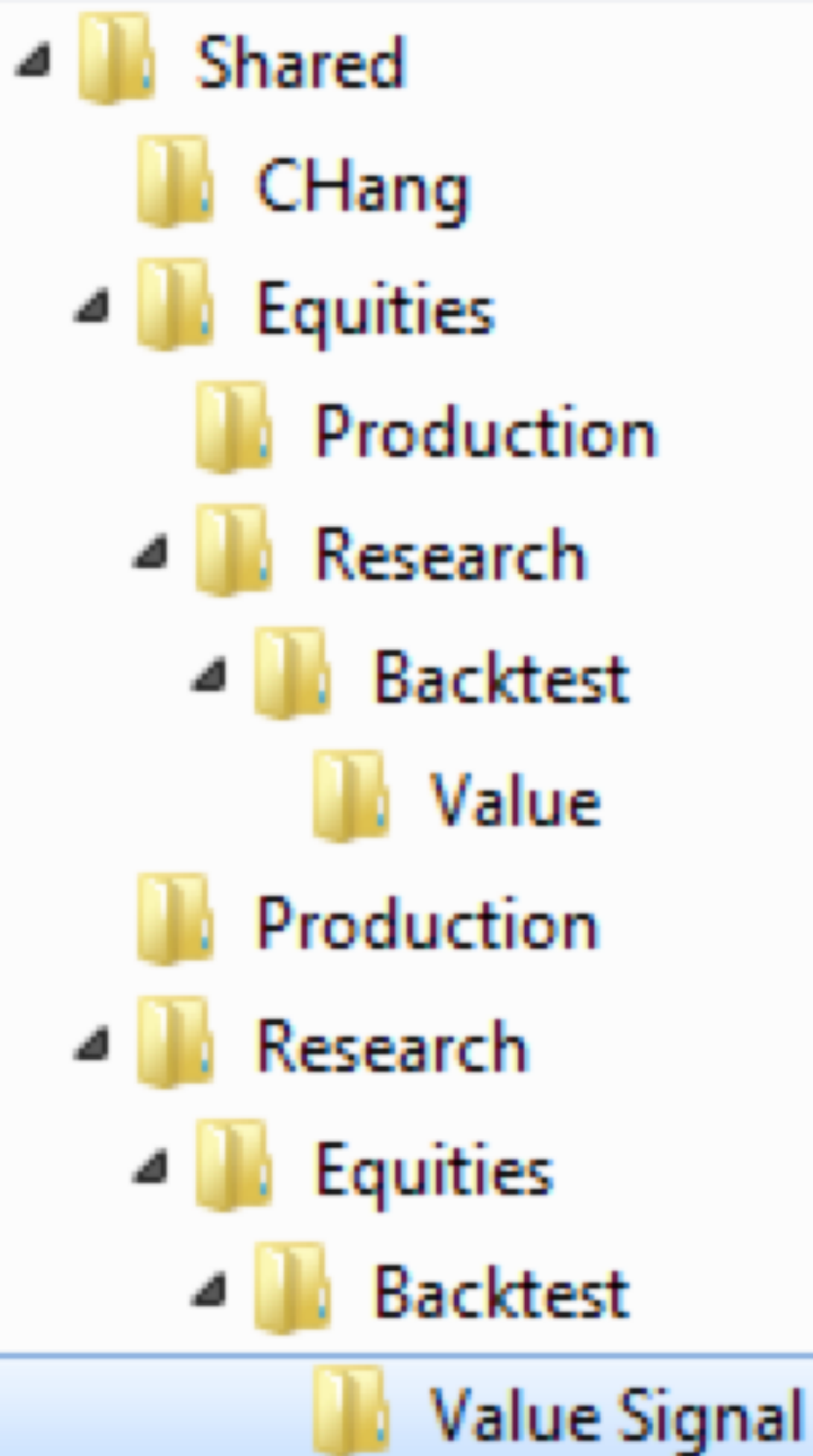


Quantitative finance (by critics)





Courtesy of: Drew Conway



How NOT to organize

- Which folder has the backtest results for the value signal?
- Friends don't let friends organize files like this

Why use the shared network drive?

- Backup
- Collaboration
- Presentation
- Version control

Solutions

- Right tools
- Save code with results
- Save version numbers for key libraries

O'REILLY

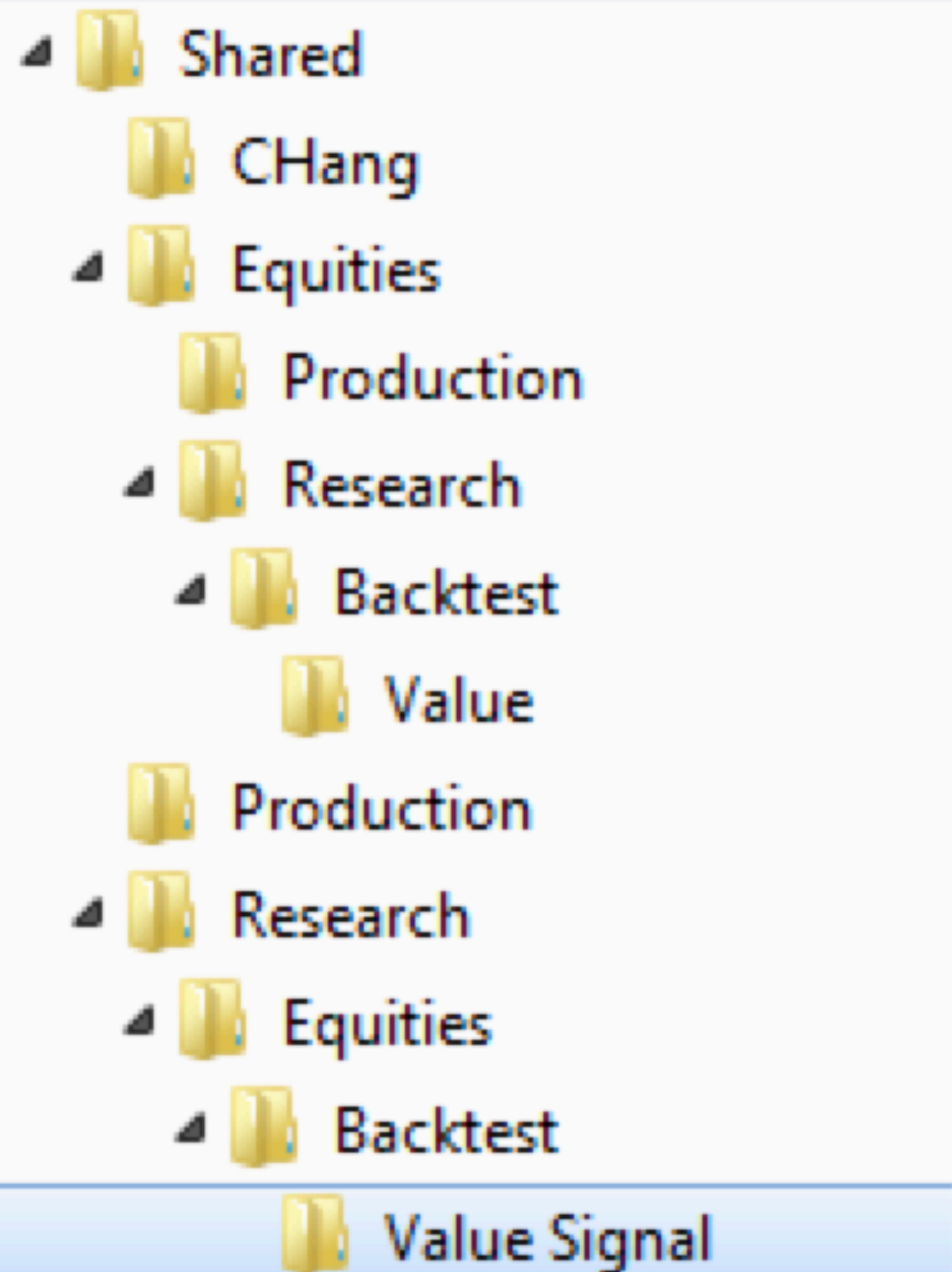
Strata
CONFERENCE

+

HADOOP
WORLD

What are “the right tools”?

- Backup and sync
- Easy sharing
- Version control



Name

backtest_results_orig.xlsx

backtest_results_back.xlsx

backtest_results.xlsx

backtest_results - Copy.xlsx

backtest_results_20121023.xlsx

backtest_results.xlsx.bak

backtest_results_changshe.xlsx

backtest_results_CS_20120915.xlsx

backtest_results.xls

backtest_results_new.xlsx

O'REILLY

Strata
CONFERENCE

+ HADOOP
WORLD

(Version) Control all the things!

- Code
- Configurations
- Data

OREILLY

Strata
CONFERENCE

+

HADOOP
WORLD

Code version control

- Choose the right tool for your workflow
- High activation energy for financial researchers



O'REILLY

Strata
CONFERENCE

+

HADOOP
WORLD

Code version control

- Stop building models in Excel
 - Traceability
 - VBA, or “I don’t want to live on this planet anymore”

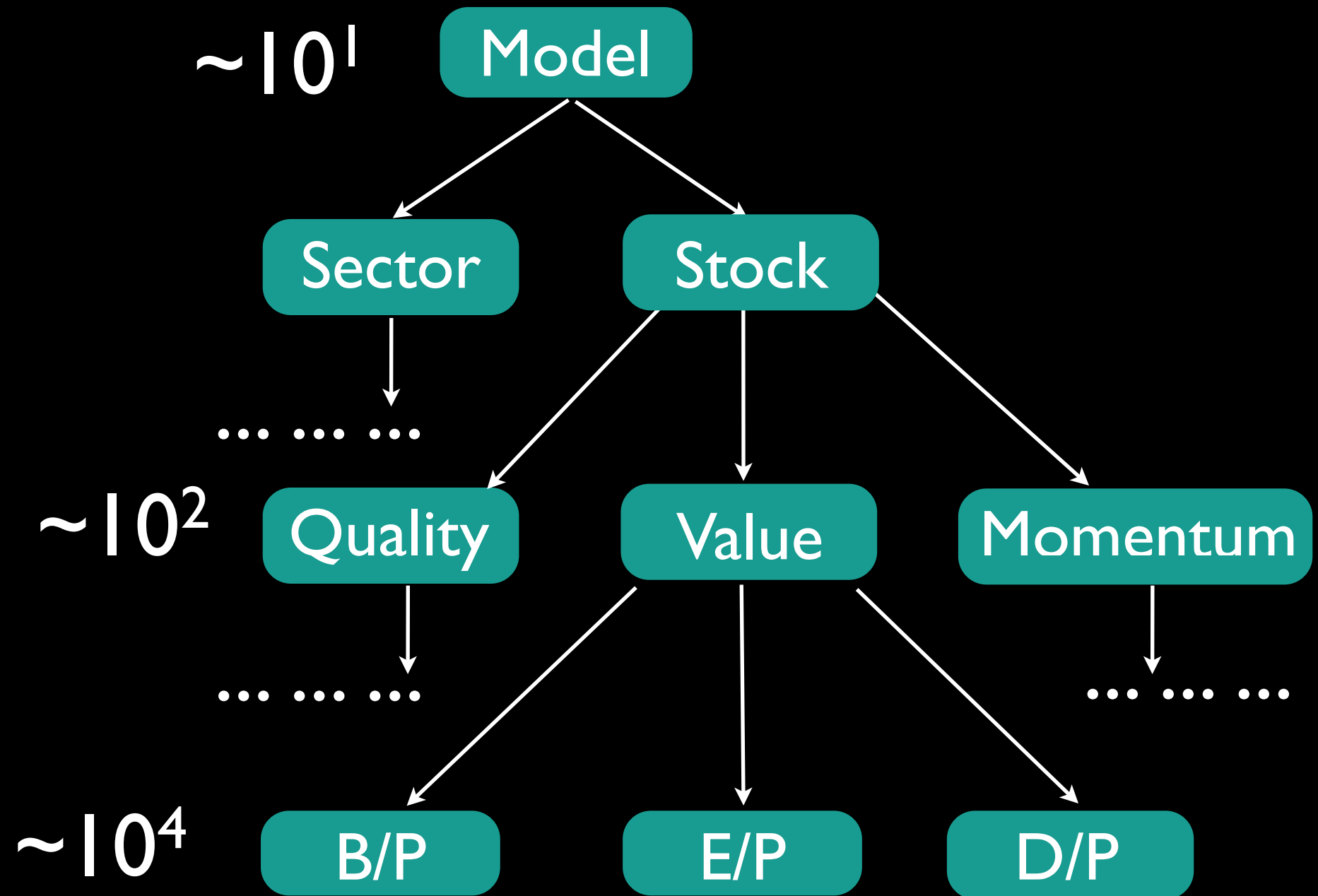
Code versioning tools

- Distributed vs Centralized
- How do you collaborate?
- What's least disruptive to your workflow?



~~parameter hell~~ Configurations

- Factor tree
- Model settings
- Factor settings
- Risk model
- Trading cost model
- Constraints



Configuration version control

- Too many dials
- Version changes often unrecorded
- Multiple time dimensions



call
database

```
2 def get_factor_config(factor, date):
3     sqlstr = """\
4         SELECT * FROM Factors
5         WHERE FactorName='%s' AND
6             '%s' BETWEEN StartDate and EndDate
7         """
8     rs = database_query(sqlstr % (factor, date), CON)
9     return FactorConfig.from_sql(rs)
10
11
12 class Factor(object):
13
14
15     def config(self, date):
16         value = get_factor_config(self.id, date)
17         return value
```

wrapper

O'REILLY

Strata
CONFERENCE

+

HADOOP
WORLD

```
2 def get_factor_config(factor, date):
3     sqlstr = """\
4         SELECT * FROM Factors
5         WHERE FactorName='%s' AND
6             '%s' BETWEEN StartDate and EndDate
7         """
8     rs = database_query(sqlstr % (factor, date), CON)
9     return FactorConfig.from_sql(rs)
10
11
12 class Factor(object):
13
14     _config = None
15
16     def config(self, date):
17         if (self._config is None or
18             self._config.StartDate > date or
19             self._config.EndDate <= date):
20             self._config = get_factor_config(self.id, date)
21         return self._config
```

memoize

O'REILLY

Strata
CONFERENCE

+

HADOOP
WORLD

Configuration versioning

- Part of the problem can be solved with additional date information
- But lots of ad-hoc wrapper code. (“Wait...which tables have start/end dates again?”)
- What if a particular research study requires changing historical values?

Designing parameter version control

- Centralized entry point for configuration querying
- Get and set version date / tag / hash
- Runs in different anchoring modes
 - Single anchor - use current configurations
 - Multiple anchor - switch at pre-defined dates
 - Floating anchor - match data date / version

**IF DATA CHANGES AND NO ONE
KNOWS ABOUT IT**



**DID IT REALLY
CHANGE?**

memegenerator.net

Data

- Scale problems
- Boundaries of control
- Proprietary point-in-time databases

O'REILLY

Strata
CONFERENCE

+

HADOOP
WORLD

Data versioning

- Distribution
- Diff / Merge
- Performance

Data versioning design issues

- Store locally only as needed
- Need a structure aware diff tool
- Scale of data means cannot all be “full” versions
- Quick access means cannot all be change-sets

Data versioning design issues

- Need some full save-points (e.g., per day or month)
- Just store change-sets for other points
- A way to “upgrade” partial save-points to full if accessed very frequently

Dependency hell

- Firm-wide research platform
 - Must have full scientific computing stack
 - Cloneable VMs
- New research tools: <http://www.pgbovine.net/cde.html>



Testing in the financial industry

- Cavalier attitude
- Lack of good processes
- Over-reliance on compile time checks
- Reinventing the wheel, on purpose!

Testing

- Testing is worth the time
- Versioning \Leftrightarrow Testing
- Continuous integration

O'REILLY

Strata
CONFERENCE

+

HADOOP
WORLD

Unit testing

- Independent of configuration and data versions
- Data loading/cleaning/munging code
- Data transformation/computation code

Loading/cleaning/munging

```
assert data.name == expected
```

```
assert (isnull(data) == exp).all()
```

```
assert data.shape == expected
```

```
assert I am using Python + pandas
```

O'REILLY

Strata
CONFERENCE

+

HADOOP
WORLD

Core computations

```
assert result.std() == expected
```

```
assert calc() == alternate(ddof=1)
```

```
assert ar_reg(test_data) == exp
```

Data testing

```
assert has_dataitems(expected)
```

```
assert data.dtype == expected
```

```
assert data.count() == expected
```

Data testing

```
assert abs(returns) < expected  
  
assert sector_code in GIC_CODES  
  
assert problem_data == expected
```

Model testing

- `assert model(test) == expected`
- `assert model_today == model_yest`
 - Data versioning
- Reasonable run time

Conclusion

- Organize, with an eye for collaboration
- Version, EVERYTHING
- Test, fully and rigorously

O'REILLY

Strata
CONFERENCE

+

HADOOP
WORLD