

MODELING  
MUDDY  
DATA

IT'S ALL ABOUT ME

ME ME ME ME ME ME



<http://www.flickr.com/photos/bixentro/2091438688/>





<http://www.flickr.com/photos/dasqfamily/348368964/>

3

Before diving into the nasty, muddy, non-relational, confusing details it's important to understand relational theory and how we use it to model data.



FOR  
NORMALIZE'D



# All Attributes are atomic



<http://www.flickr.com/photos/-cavin-/2313239884/>

## Columns contain one value

6

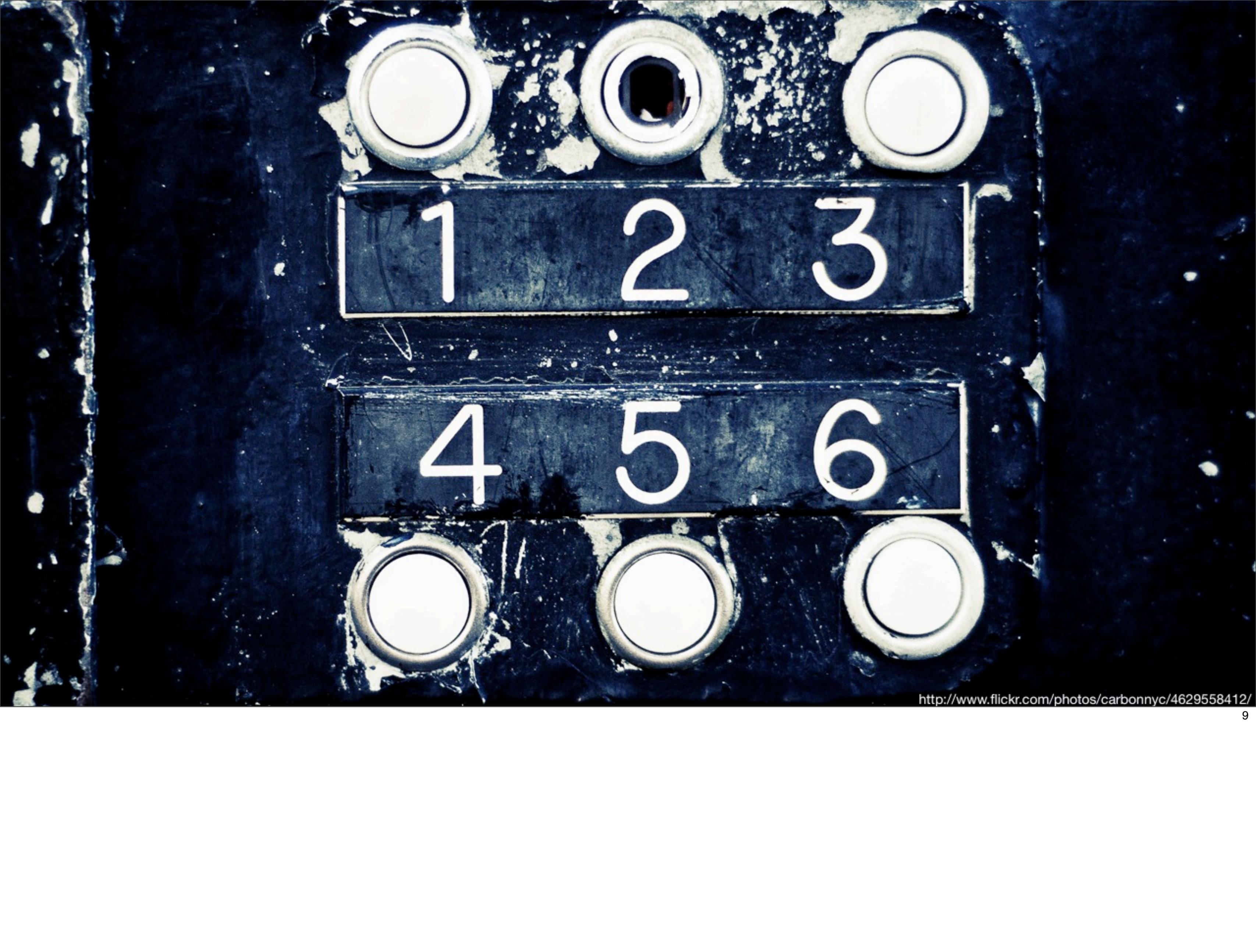
This is tricky: what is a value? Is XML a singular value if it's always consumed as a chunk? What about a binary representation of an object?

ALL ROWS  
MUST BE





Every row has the same number of columns



1 2 3

4 5 6



جیسون  
لارڈ  
فوجی  
کروز پرنسپلی



<http://www.flickr.com/photos/hokkey/324869104/>

columns  
only describe  
can  
the key, not  
other columns

columns  
only describe  
can  
the key, not  
other columns





[http://www.flickr.com/photos/x-ray\\_delta\\_one/4265173382/](http://www.flickr.com/photos/x-ray_delta_one/4265173382/)



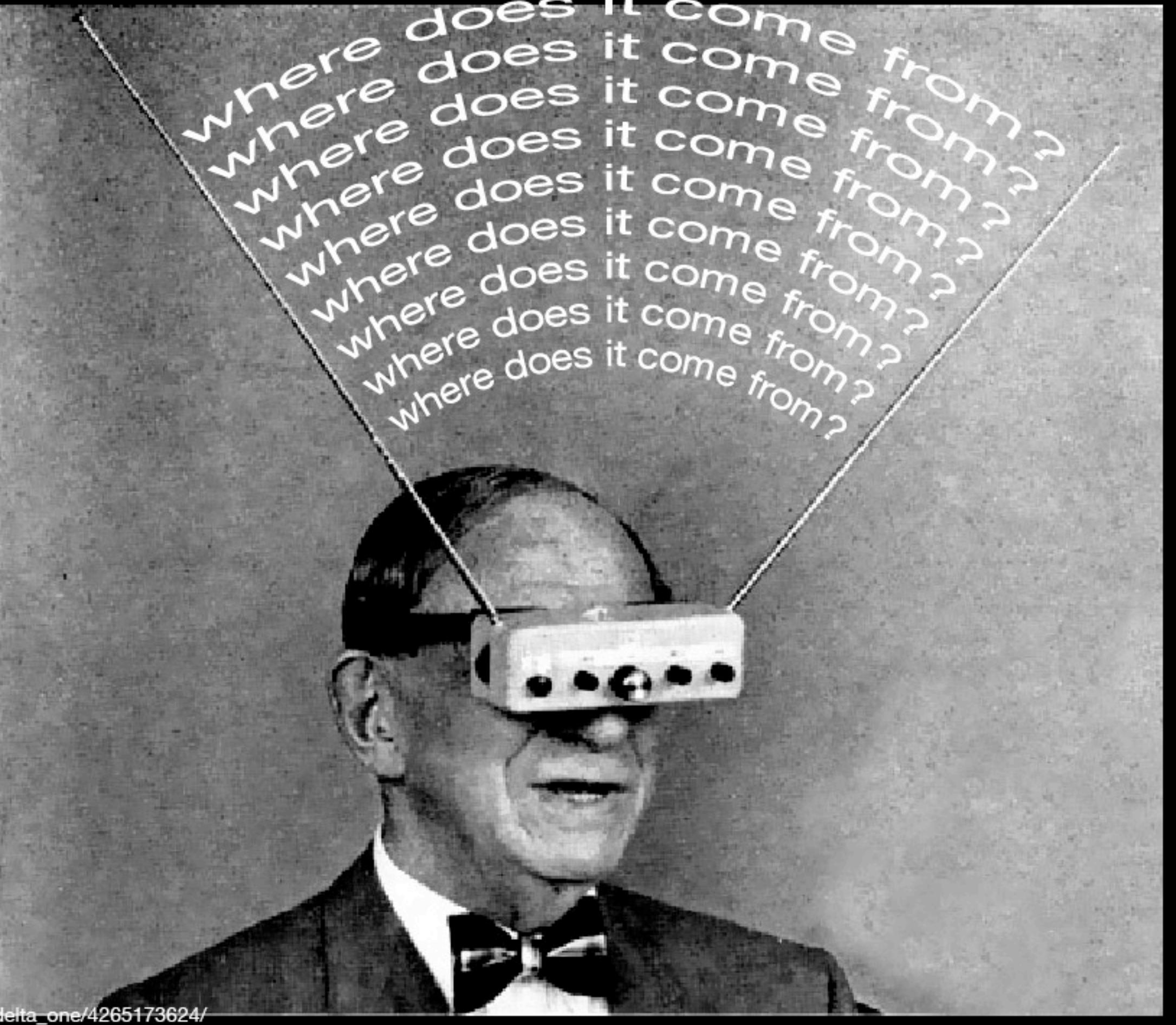
Muddy data.  
What is it?



<http://www.flickr.com/photos/dariusdunlap/749610787/>

16

It's data that doesn't fit  
Something that isn't easily stored and retrieved from 3NF  
Complex attributes  
You might call this flexible/dynamic data



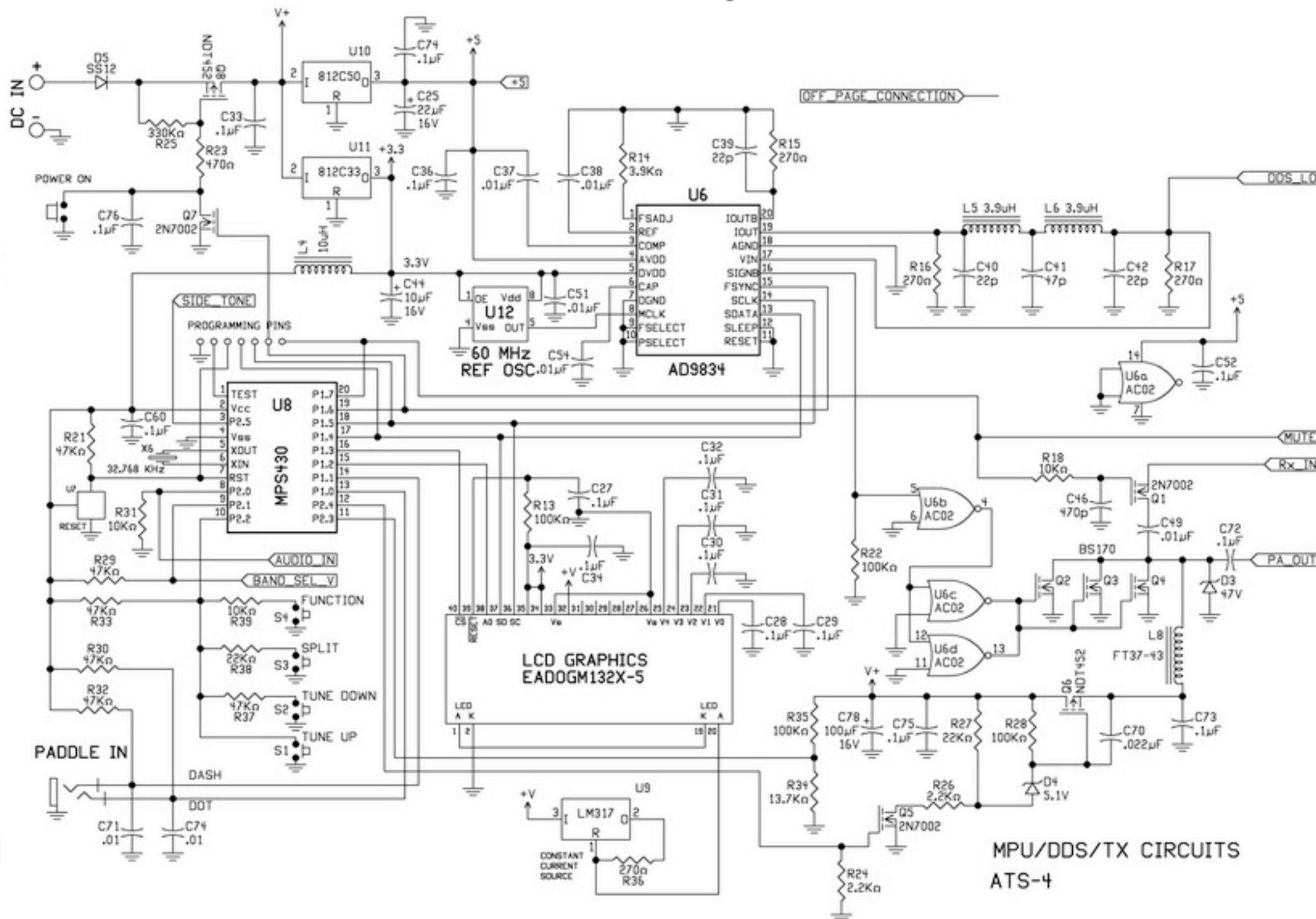
[http://www.flickr.com/photos/x-ray\\_delta\\_one/4265173624/](http://www.flickr.com/photos/x-ray_delta_one/4265173624/)

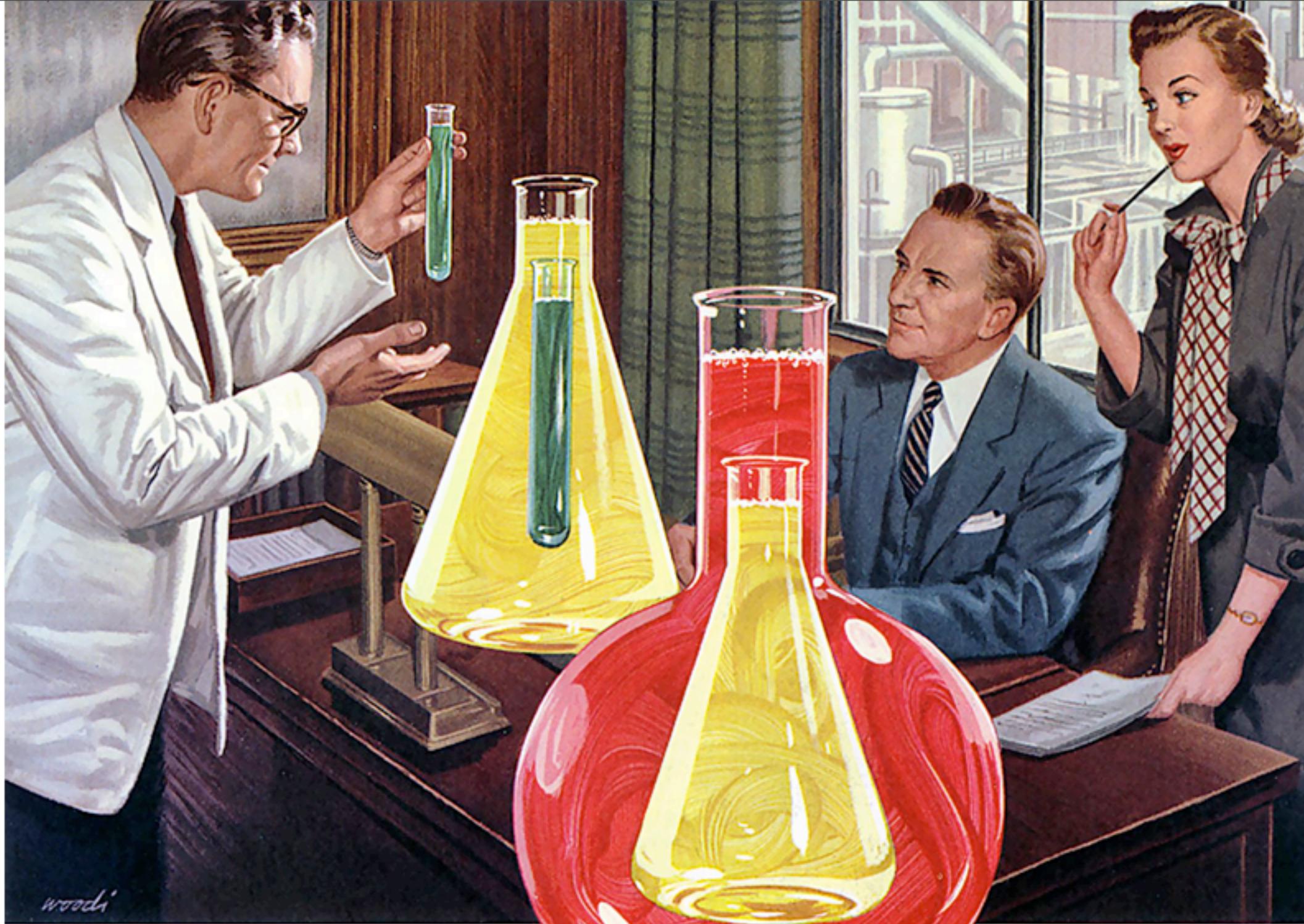
17

Ask questions!

User sessions  
Custom data

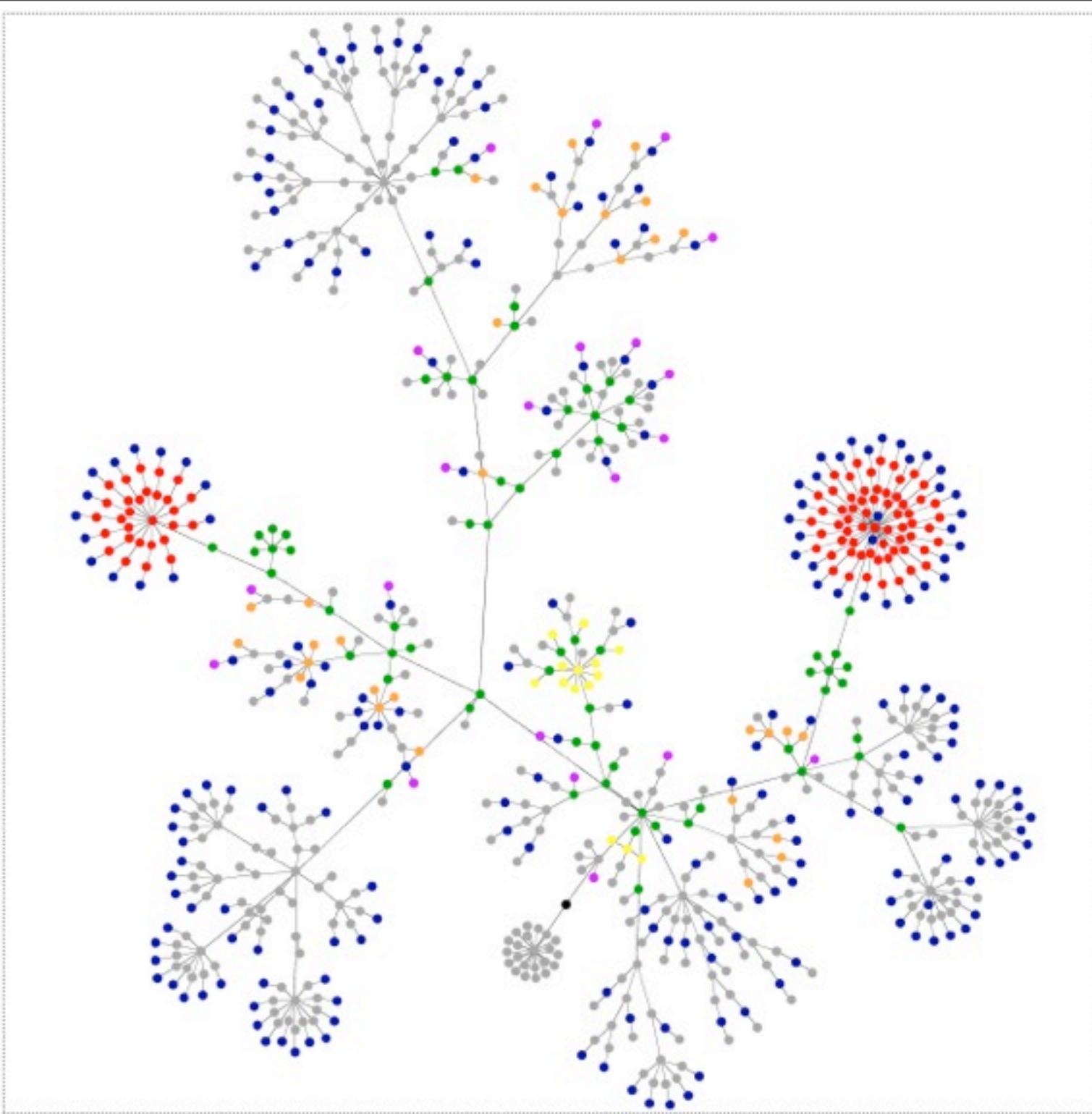
# How would you do it?





# common solutions

[http://www.flickr.com/photos/x-ray\\_delta\\_one/3968092988/](http://www.flickr.com/photos/x-ray_delta_one/3968092988/)



# entity:attribute:value

<http://www.flickr.com/photos/thefangmonster/352461415/>

20

tracks relevant facts – only record exactly what you need  
use this with many individual attributes  
requires many descriptor tables to be useable  
the physical schema is radically different from the logical schema – impedance mismatch  
business logic is in metadata instead of database schema  
when to use – sparse attributes, highly dynamic schemas

take a look at VistA database from the Veterans Health Administration for a real world example – lest you think this is a horrible database, it's won a number of awards and is being implemented in other hospitals

# Polymorphic Associations

e.g. "Tagging"

<http://www.flickr.com/photos/ceonyc/162586965/>

21

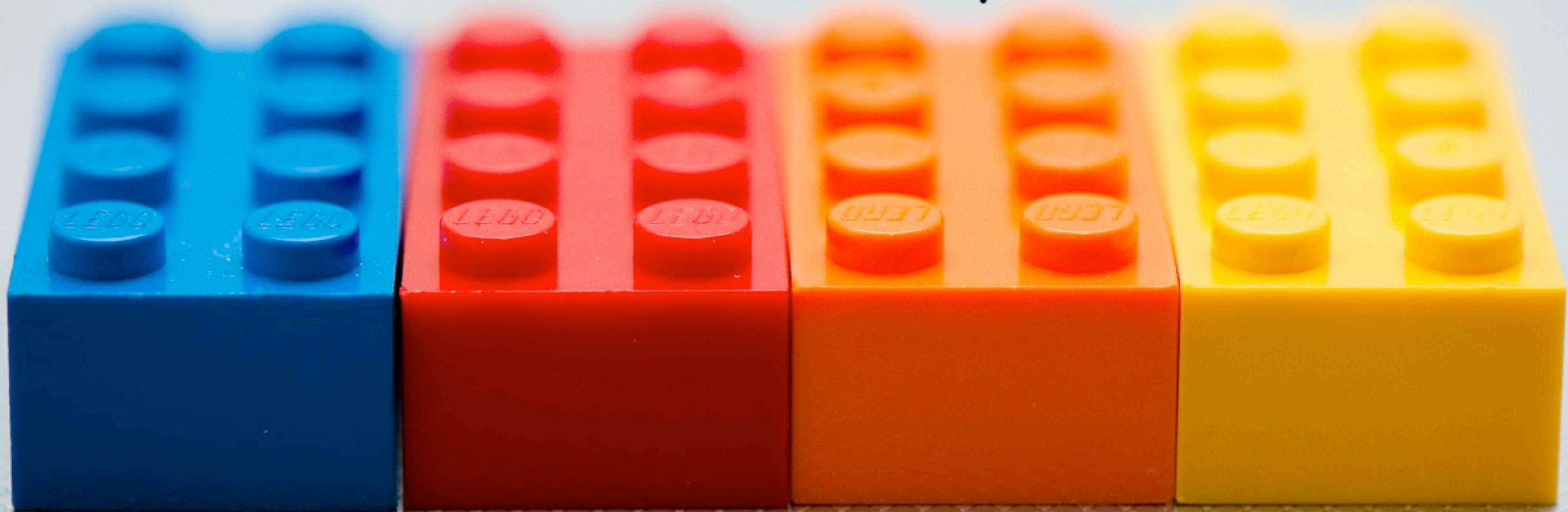
Polymorphic Associations  
somewhat specific to OO languages  
Think of tags – anything can be tagged. Thus, tags need to be generic

Let's talk about NULL.

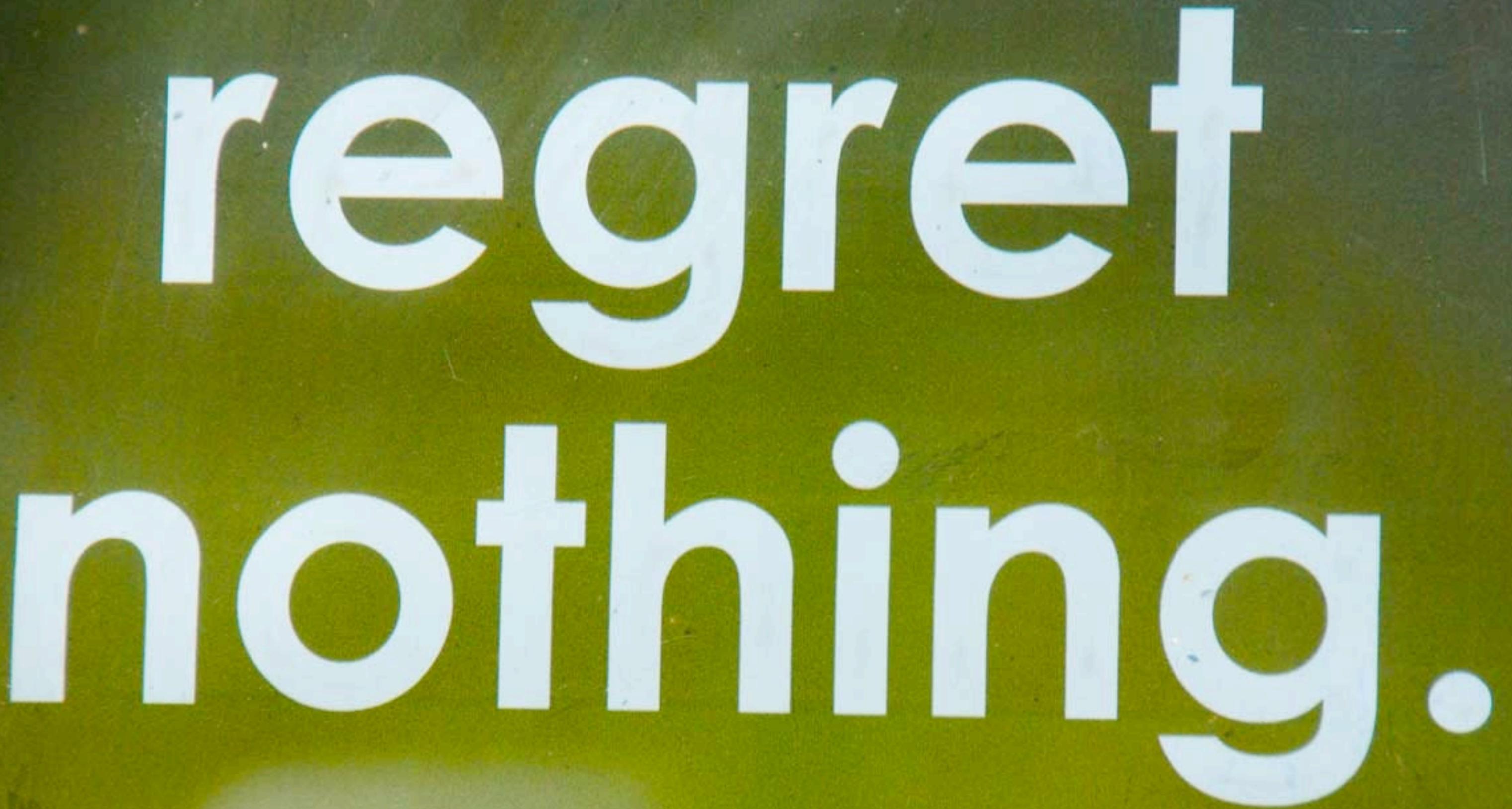
NULL is an unknown value.

You could argue that optional values are just unknown, right?

It's easy to extend your app



(like using legos made out of nothing)



regret  
nothing.

<http://www.flickr.com/photos/yourdon/3550794139/>

24

If NULLs are great, why would we not want to use them?

The purpose of a database is to keep track of data – things we know about. Why bother with things we don't know about?

Becomes difficult to distinguish between NULL and "set to nothing" and we end up with "No data entered" values.

Code becomes complex with checks for NULL values.

Query plans, as a result, become more complex

Three value logic is a hoot! (TRUE, FALSE, and UNKNOWN)

Wide tables are also a PITA to understand, index, and deal with.

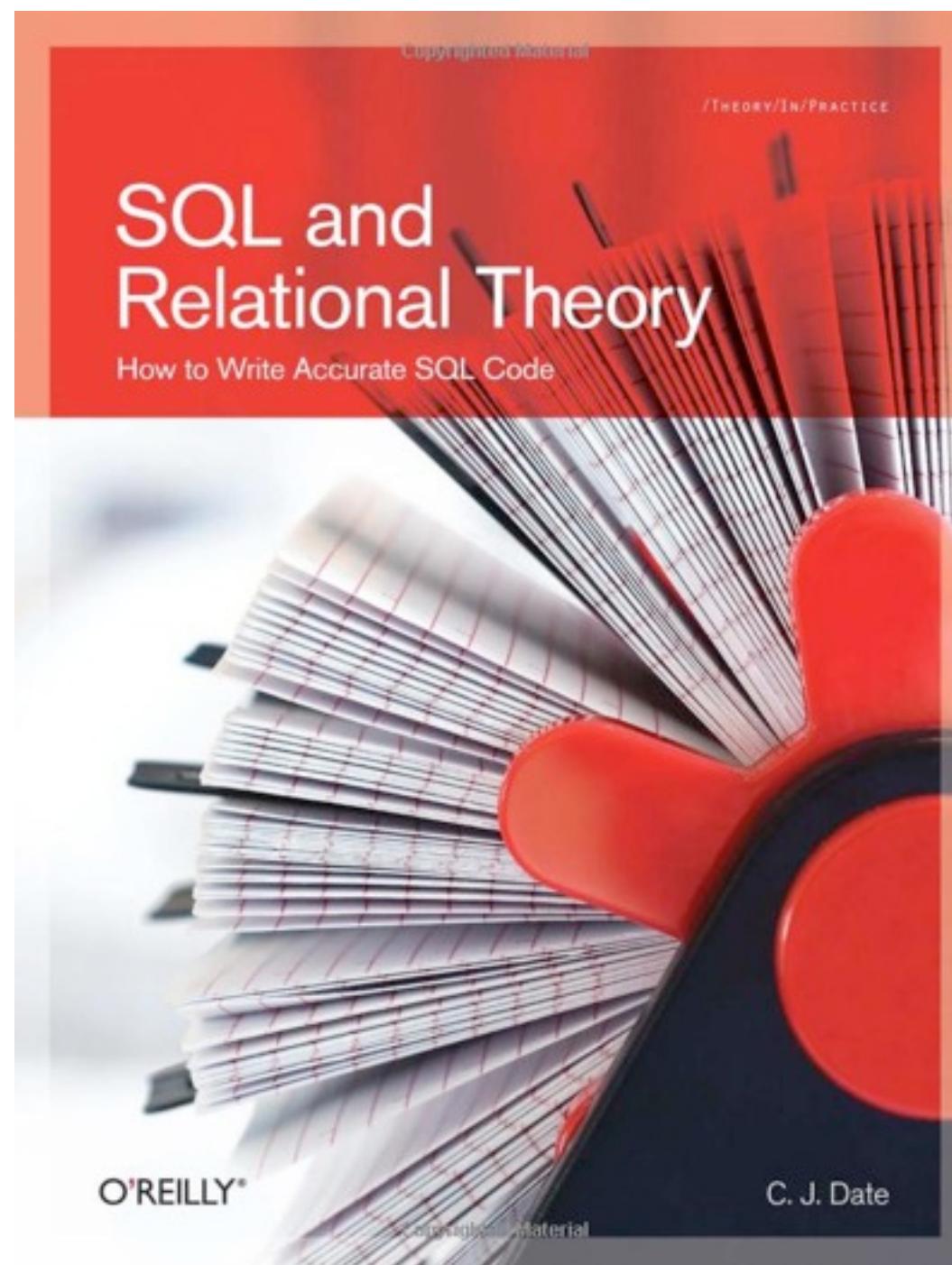
# The Academic Reason

Let  $T_1, T_2, \dots, T_n$  ( $n \geq 0$ ) be type names, not necessarily all distinct. Associate with each  $T_i$  a distinct attribute name,  $A_i$ ; each of the  $n$  attribute-name/type-name combinations that results is an attribute. Associate with each attribute an *attribute value*  $v_i$  of type  $T_i$ ; each of the  $n$  attribute/value combinations that results is a component. The set of all  $n$  components thus defined,  $t$  say, is a tuple value (or just a tuple for short) over the attributes  $A_1, A_2, \dots, A_n$ . The value  $n$  is the degree of  $t$ ; a tuple of degree one is unary, a tuple of degree two is binary, a tuple of degree three is ternary, ..., and more generally a tuple of degree  $n$  is  $n$ -ary. The set of all  $n$  attributes is the heading of  $t$ .

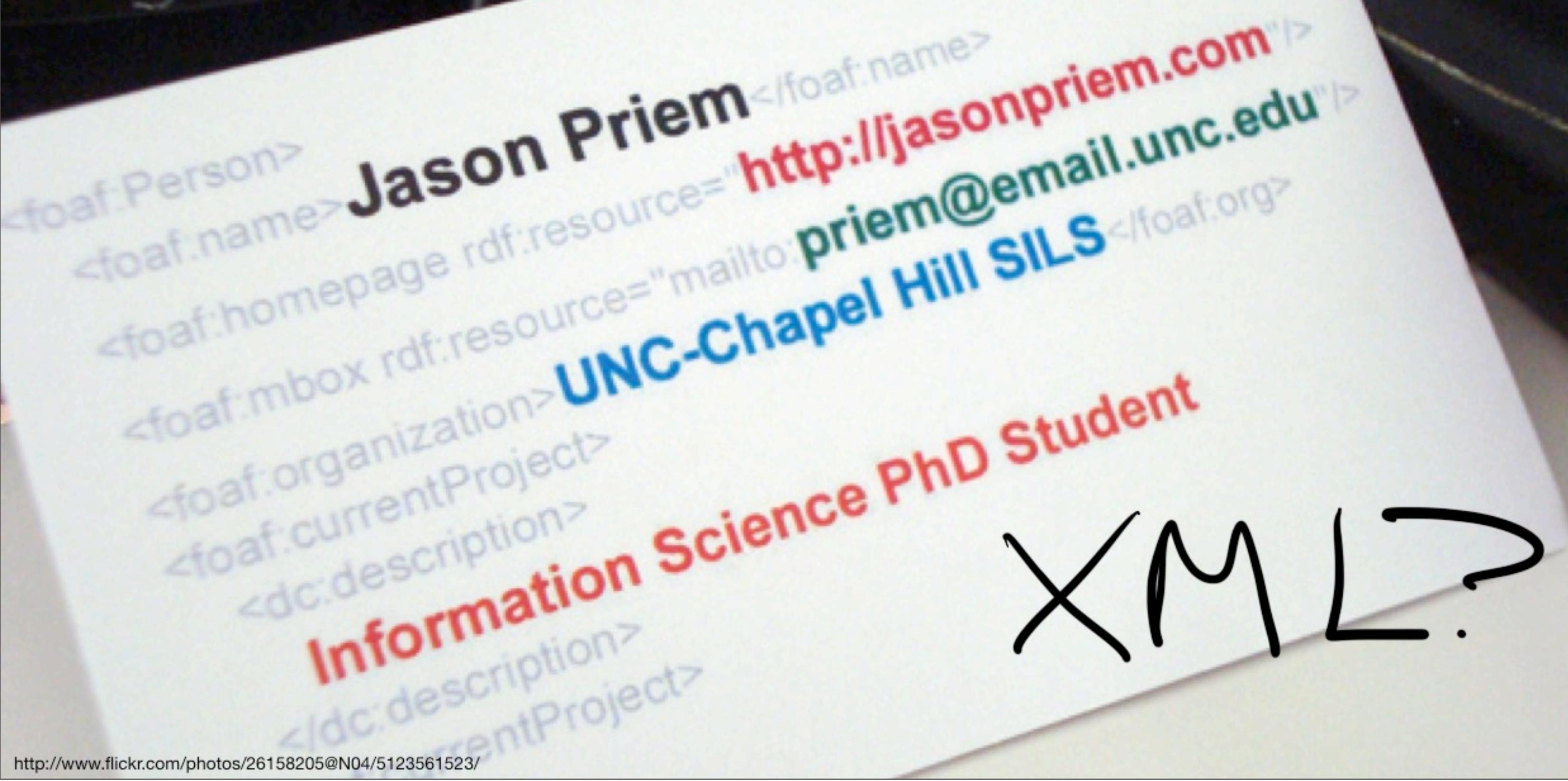
SQL and Relational Theory by C. J. Date. Copyright 2009 C. J. Date, 978-0-596-52306-0.

By definition a tuple (row) contains a value of the appropriate type for each attribute (column). If a NULL is an unknown, it's not a value (NULL != NULL)  
Thusly, a table with NULLs isn't really a table at all and is some kind of horrible degenerate monster.

# For More Hatred of NULLs...



# What about



Pros:  
Flexible, document structure  
Rules can be enforced via a schema  
Can be validated both in SQL Server and in the application tier  
Can support rich data representations relatively simply

Cons:  
Can't use native temporal datatypes in SQL Server 2005 – XML temporal data requires a timezone (you need a workaround)  
Breaks the rules of normalization (if you're using NULLs, who cares?)



28

This is sometimes called 'table per subclass'

Pros:

core data is stored in one table

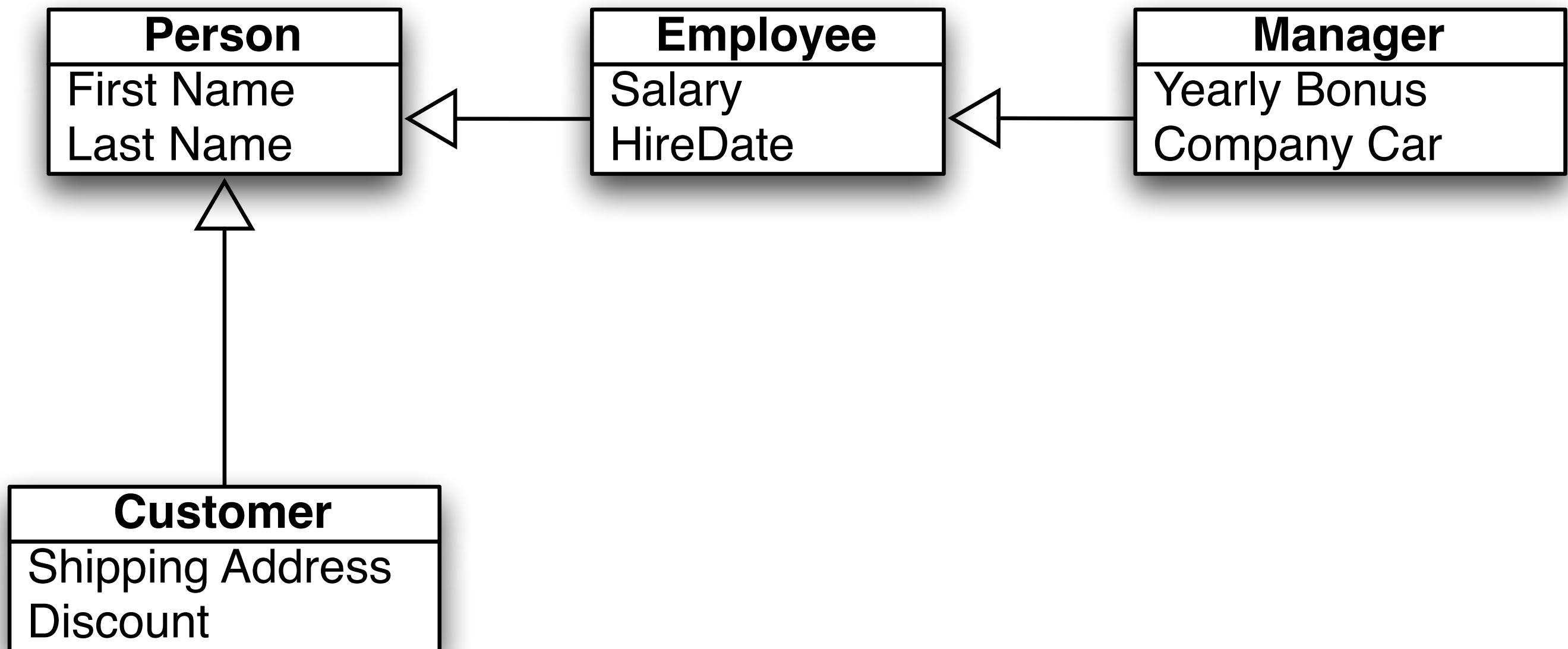
divide optional data into discrete feature groups then create a separate table for each feature group

Cons:

many 1:1 joins

complex querying

what happens when you sub-class a sub-class?



~~THERE IS~~ No



<http://www.flickr.com/photos/cubanrefugee/2820985266/>

30

Ask yourself questions about how the data will be written and how it will be queried. Once you are comfortable with normalization, usage patterns should dictate the data model, not adherence to "the rules".

Look at frequently run queries, slow queries, indexing patterns, strange/complex joins, tight coupling, cardinalities, access patterns, and shoehorning

there are other ways...



32

HBase – modelled after Google's BigTable. Fixed column "families", fuzzy columns

MongoDB – document database, has secondary indexes, very much like an RDBMS

CouchDB – document database – flexible data, indexing, stored views, multi-master replication

Redis – in memory key/value store with support for complex data types (lists and sets), new versions support replication, flexible backends

Riak – distributed key/value store with flexible storage backends, distributed, replicated

So why go down this dark and mysterious route? All of these options provide you with varying degrees of flexibility, querying, and scalability. Scaling EAV is hard. Scaling XML querying is hard. Don't be afraid to do something insane.

# IMPORTANT INFORMATION

This work is licensed under a  
Creative Commons Attribution-ShareAlike 3.0 Unported License.



This means that you can make it better, but  
you have to keep sharing.

<http://github.com/peschkaj/Muddy-Data>

<http://facility9.com>

You should follow me on twitter: <http://twitter.com/peschkaj>