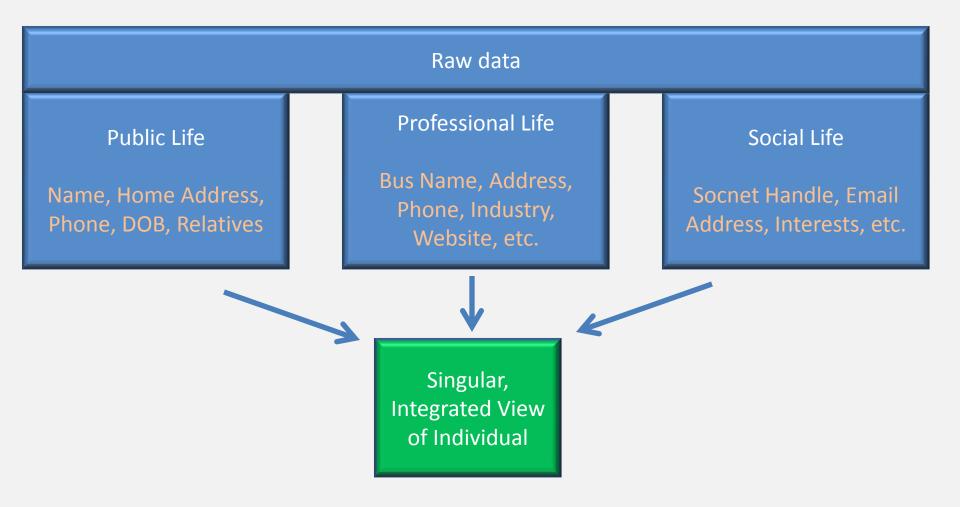
Navigation of Cascading Data Errors

Kevin Marcus, CTO, Intelius

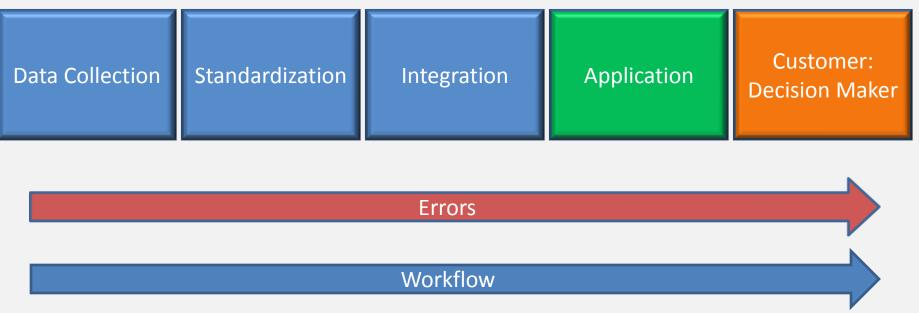
The Problem



- Apps want a single id: "People, Business, Assets"
- No universal globally available identifiers
- Data scattered in places, formats, attributes

Workflow

- Output of each step is input to next.
- Impact of errors are magnified as they progress.
- Any manipulation of original data has potential to introduce error



Data Collection

- Easy to do, Hard to do well.
- Entry Errors
- Digitization Errors
- False Information
- Precision
- Verification can be difficult.

FirstName	MiddleName	Las
ROBERT	М	BE
BOB	S	BE
ROBERT	MARTIN	BE
ROBERT	М	BE
BOB	SHELTON	BE
ROBERT	М	BE
ROBERT	S	BE
ROBERT	М	BE
ROBERT	М	BE

MZ	R	JANE	BE
MA	RY	JEAN	ZE
54	6	194	10615
54	4	198	30309
54	2	194	10600

OLD WOODINVILLE DUVALL RD NE OLD WOODINVILLE DUVALL R OLD WOODINVILLE DUVA RD NE OLD WOODINVILLE DUVAL

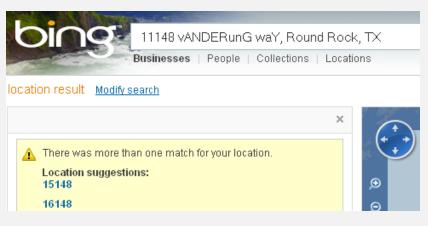
Businesses may not be incented to collect high quality data!

Standardization

- Formatting ambiguities
- Some industries more advanced than others (Finance, HR)
- No "global data model"
 - Best friend: Multisampling convergence over time.
- Address standardization is hard and still unsolved.

Standardize Address: "11148 vANDERunG waY, round rock, tx"





Search Result

Your Input Search Result

11148 vANDERunG waY 11148 vANDERunG waY round rock TX

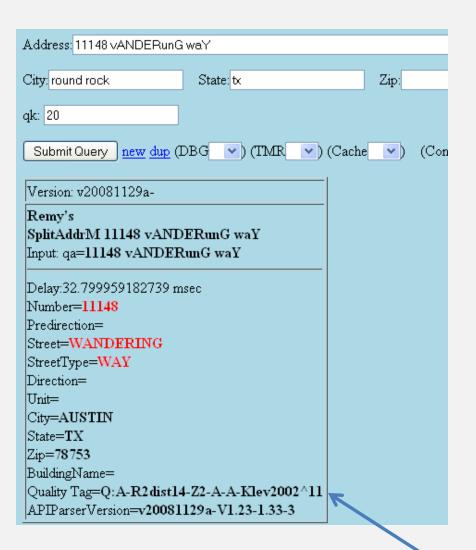
Search warnings

4.1 Address not found 6.1 Multiple streets match





Best Guess? 20 miles away – 11148 Wandering Way





"Low quality match"

Record Integration is Hard

- Conflation
 - Are two records the same?
 - Statistical models are not perfect.
- De-dupe: Apps want a single ID!
- Orphans and Black Holes
 - Orphan: record unlinkable to others.
 - Black Hole: record matches too many others.
- Data Symmetry
 - Data interface restrictions. (e.g. CLID)
 - Offset with better conflation, caching

```
city': 'SomeCity',
                         'postal code': '12345',
                         'state': 'AA'}],
     'all tags': [ 'ABC county',
                     'member of ABC County Association Of Realtors',
                     'works at ABC Inc',
                     'realtor'.
                     'AA',
                     'real estate agent',
                     'ABC']
      'emails': [{'value': 'email@email.com'}],
      'names': [ { 'display name': 'L B',
                    'first name': 'L',
                     'last name': 'B',
      'phone numbers': [{'number type': 'fax', 'sources': [0], 'value': '1234567890'}],
      'tags': [ {'value': 'ABC county'},
                {'value': 'member of ABC County Association Of Realtors'},
                 {'value': 'works at ABC Inc'},
                 {'value': 'realtor'},
                 {'value': 'real estate agent'}]}
   'addresses': [ { 'address street1': '123 Street',
                      'city': 'SomeCity',
                      'postal code': '12345',
                      'state': 'AA'}],
   'all tags': ['Physician Assistant', 'General practitioner', 'ABC', 'physician Assistant', 'AA'],
   'jobs': [{'job title': 'Physician Assistant'}],
   'names': [{'first name': 'L', 'last name': 'B', 'niddle name': 'I'}],
   'tags': [ {'value': 'physician Assistant'},
              {'value': 'General practitioner'}]}
              FirstName MiddleName LastName
PartyID
                                                     HouseNum StreetName City
                                                                                    State
                                                                                            Zip
            22 L
                                     В
                                                             456 Street
                                                                            SomeCity AA
                                                                                                12345
```

{ 'addresses': [{ 'address street1': '0 Street'

22 L

96 L

SomeCity AA

SomeCity AA

123 Street

0 Street

12345

12345

В

В

Applications and Decisions

- Historical: Data + App often originated at same organization.
 - No longer true!
 - Data is bought, sold, processed and sliced from various sources!
- Apps: Guilty until proven innocent
- Computers should help a human decide
 - Should not be the final judgment!
 - Al not quite good enough yet ©

M W M II: Parental Control!

- Unique name: Only 3 people in US; 2 in WA!
- 43 Criminal records
- 10 Civil cases
- DOB available on many not all!
- "Use suffix when writing signature"
- Obvious issues: "Given the filing date, would it make sense for II?"





Future Needs

- More data collection transparency.
- Better error detection and correction tools.
- Improved record conflation mechanisms.
- International influences
 - Privacy and collection issues vary
 - Schemas may need to change.
- Easier global "fix" for individuals affected.
 - Centralized repository?
 - Better authentication and validation systems?