



## Genetic data

```
(raccoon:19.19959,bear:6.80041):0.84600;((sea lion:11.99700,seal:12.00300):7.52973;((monkey:100.85930,cat:47.14069):20.59201,weasel:18.87955):2.09460):3.87382,dog:25.46154);(Bovine:0.69395,(Gibbon:0.36079,(Orang:0.33636,(Gorilla:0.17187,(Chimp:0.19268,Human:0.11927):0.08386):0.06124):0.115057):0.54939,Mouse:1.21460):0.10;Bovine:0.69395,(Hylobates:0.36079,(Pongo:0.33636,(G_Cori11a:0.17147,(P_paniscus:0.19268,H_sapiens:0.11927):0.08386):0.06124):0.115057):0.54939,Rodent:1.21460);
```

## Haskell HI files

```
00000000: 0001 face 0000 0073 0400 0000 3600 0000 .....s...S...
00000010: 3000 0000 3500 0000 3000 0000 0000 0000 0...S...D...
00000020: 0001 0000 0000 0100 0000 0043 0001 0000 .....C...
00000030: 0002 0200 0000 0200 0000 0300 0000 0200 .....H...
00000040: 0000 0400 0000 4800 0100 0000 0200 0000 .....H...
00000050: 0502 0000 0000 0006 0000 0000 0007 0000 .....H...
00000060: 0001 0000 0000 6800 0000 0000 006f 0000 .....H...
00000070: 0000 0100 0000 0800 0000 0968 6173 6b65 .....Haske
00000080: 6e6c 3938 0000 0007 4350 5554 696d 6500 1198....CPUTime.
00000090: 0000 0462 6173 6500 0000 0847 4843 2e42 ...base...GHC.B
000000a0: 6173 6500 0000 0e47 4843 2e46 6f72 6569 ase...GHC.Forei
000000b0: 676e 5074 7200 0000 0e53 7973 7465 6d2e gnPtr...System.
000000c0: 4350 5554 696d 6500 0000 0a67 6574 4350 CPUTime...getCP
000000d0: 5554 696d 6500 0000 1063 7075 5469 6d65 UTime...cpuTime
000000e0: 5072 6563 6973 696f 6e Precision
```

## Ad hoc data from AT&T

Name & Use	Representation	Size
<b>Web server logs (CLF):</b> Measure web workloads	Fixed-column ASCII records	≤ 12 GB/week
<b>Sirius data:</b> Monitor service activation	Variable-width ASCII records	2.2GB/week
<b>Call detail:</b> Detect fraud	Fixed-width binary records	~7GB/day
<b>Altair data:</b> Track billing process	Various Cobol data formats	~4000 files/day
<b>Regulus data:</b> Monitor IP network	ASCII	≥ 15 sources, ~15 GB/day
<b>Netflow:</b> Monitor IP network	Data-dependent number of fixed-width binary records	>1Gigabit/second

## And many others...

- Gene ontology data
- Call detail data
- Cosmology data
- Netflow packets
- Financial trading data
- DNS packets
- Telecom billing data
- Java JAR files
- Router config files
- Jazz recording info
- ...
- System logs

## Technical challenges

- Data arrives "as is."
- Documentation is often out-of-date or nonexistent.
  - Hijacked fields.
  - Undocumented "missing value" representations.
- Data is buggy.
  - Missing data, human error, malfunctioning machines, race conditions on log entries, "extra" data, ...
  - Processing must detect *relevant* errors and respond in *application-specific* ways.
  - Errors are sometimes the *most* interesting portion of the data.
- Data sources often have high volume.
  - Data may not fit into main memory.

## Existing approaches

- Lex/Yacc
  - No one we have encountered uses them for ad hoc data.
- Perl/C
  - Code brittle with respect to changes in input format.
  - Analysis often ends up interwoven with parsing, precluding reuse.
  - Error code, if written, swamps main-line computation. If not written, errors can corrupt "good" data.
  - Everything has to be coded by hand.
- Data description languages (PacketTypes, Datascript)
  - Binary data

## types to the rescue!

Relational and XML data are relatively easy to manage (partly) because schema exist to describe the data.

Relational Data	Relational Schema
XML	XML Schema
Ad Hoc Data	<i>Physical Types</i>

**Thesis:** Types can facilitate ad hoc data management, and the types developed for in-memory values are suited to the task.

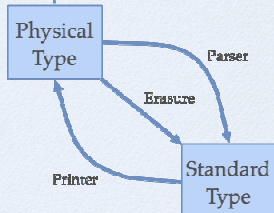
## Our approach: PADS

- Data expert writes declarative description of data source:
  - Physical format information
  - Semantic constraints
- Many data consumers use description and generated parser.
- Description serves as living documentation.
- Parser exhaustively detects errors without cluttering user code.
- From declarative specification, we generate auxiliary tools.

## Typing ad hoc data

```
"CREAJEUS",157913,-1.84823,Quadrillion Stu,4
"CREAJEUS",157813,-0.33639,Quadrillion Stu,4
"CREAJEUS",157913,-1.64932,Quadrillion Stu,4
"CREAJEUS",188213, 1.8537,Quadrillion Stu,4
```

Described by



## Base types

```
"CREAJEUS",197313,-0.456483,Quadrillion Stu,4
"CREAJEUS",197413,-0.482265,Quadrillion Stu,4
```

String, Int, Float

## Tuple types

```
"CREAJEUS",197313,-0.456483,Quadrillion Stu,4
"CREAJEUS",197413,-0.482265,Quadrillion Stu,4
```

String \* Int \* Float \* String \* Int

## Singleton types

```
"CREAJEUS",197313,-0.456483,Quadrillion Stu,4
"CREAJEUS",197413,-0.482265,Quadrillion Stu,4
```

```

\'\' * String * \'\' * \',\'
 * Int * \',\'
 * Float * \',\'
 * String * \',\'
 * Int

```

Where we write \',\' for S(\',\').

## Simple dependent types

```

"TEA,BUS",197313,-0.482285,Quadrillion Btu,4
"TEA,BUS",197413,-0.482285,Quadrillion Btu,4
  
```

```

\\"" * String(\\"" ) * \\"" * \',
  * Int                      * \',
  * Float                    * \',
  * String(',,')              * \',
  * Int
  
```

## Records

```

"TEA,BUS",197313,-0.482285,Quadrillion Btu,4
"TEA,BUS",197413,-0.482285,Quadrillion Btu,4
  
```

```

{
  source: String(\\""),
  date: Int,
  measurement: Float,
  units: String(',,')
  order: Int
}
  
```

## Unions

```

Southern California Regional Railroad Authority,"Los Angeles, CA",
7,45,46,46,47,49,51,0,45,46,46,47,49,51
Connecticut Department of Transportation,"New Haven, CT",
U,U,U,U,U,U,8,U,U,U,U,U,U,8
Tri-County Commuter Rail Authority,"Miami, FL",
U,U,U,U,U,U,18,U,U,U,U,U,U,18
  
```

Anonymous:

```

'U' + Int
  
```

Named:

```

type OptInt = unavailable of 'U'
             | available of Int
  
```

## Arrays/Lists

```

Southern California Regional Railroad Authority,"Los Angeles, CA",
7,45,46,46,47,49,51,0,45,46,46,47,49,51
Connecticut Department of Transportation,"New Haven, CT",
U,U,U,U,U,U,8,U,U,U,U,U,U,8
Tri-County Commuter Rail Authority,"Miami, FL",
U,U,U,U,U,U,18,U,U,U,U,U,U,18
  
```

```

type OptInt = unavailable of 'U'
             | available of Int

type counts = OptInt[]$exp(',')
              term(eor)
  
```

## Dependent Types

```

scw-01.ab.ca - - [16/12/06] "GET /images/fish.gif HTTP/1.0" 200 8552
scw-01.ab.ca - DBUser [16/12/06] "GET /images/bug.gif HTTP/1.0" 200 1357
64.233.161.99 - - [16/12/26] "GET /images/plex.gif HTTP/1.0" 304 -
69.30.123.195 - - [16/12/2006] "GET /images/adjoint.gif HTTP/1.0" 304 -
  
```

```

type responseCode = { x : Int | 99 < x < 600}
  
```

## Dependent Types

```

scw-01.ab.ca - - [16/12/06] "GET /images/fish.gif HTTP/1.0" 200 8552
scw-01.ab.ca - DBUser [16/12/06] "GET /images/bug.gif HTTP/1.0" 200 1357
64.233.161.99 - - [16/12/26] "GET /images/plex.gif HTTP/1.0" 304 -
69.30.123.195 - - [16/12/2006] "GET /images/adjoint.gif HTTP/1.0" 304 -
  
```

```

type method = GET | POST | LINK | UNLINK | ...

fun check(method, major, minor) = ...

type request =
{ method : method,    '\',
  url : String('\',', " HTTP/",
  major : Int,        '\',
  minor : Int
} where check(method, major, minor)
  
```

# Value Abstraction

```
sdw-01:abcd@ - DBUser: [16/12/06] "GET /images/bug.gif HTTP/1.0" 200 1357
68.233.161.99 = [16/12/26] "GET /images/plex.gif HTTP/1.0" 304 -
Connecticut Department of Transportation , "New Haven, CT",
x,U,u,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x
```

Two representations for missing integer: 'U' and '-'. We can use value abstraction to reduce redundancy:

```
type OptInt = [x.unavailable of x
               | available of Int
type OptIntU = OptInt 'U'
type OptIntD = OptInt '-'
```

# Type Abstraction

```
sdw-01:abcd@ - DBUser: [16/12/06] "GET /images/bug.gif HTTP/1.0" 200 1357
68.233.161.99 = [16/12/26] "GET /images/plex.gif HTTP/1.0" 304 -
Connecticut Department of Transportation , "New Haven, CT",
x,U,u,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x,x
```

Two different types can be missing: Int and String. We can use type abstraction to reduce redundancy:

```
type Opt = At.lx.unavailable of x
          | available of t
type OptIntU = Opt Int 'U'
type OptIntD = Opt Int '- '
type OptStringD = λy.Opt (String y) '-'
```

# Recursive types

```
{Bovine:0.163295, (Gibbon:0.36079, (Orang:0.4636, (Gorilla:0.17147, (Chimp:0.19268, Human:0.11927):0.08386):0.06124):0.15957):0.54932, Mouse:1.21460):0.10
```

```
type Entry = {name : String('\', '\',
               dist : Float}
type Tree =
  | Leaf of Entry
  | Interior of
    ((' * Tree[] (\', '\', '\') * '\') * '\') * Float
```

# Pointers?

```
00000000: 0001 face 0000 0073 0400 0000 3600 0000 .....S.....8...
00000010: 3000 0000 3500 0000 0000 0000 0000 0000 0...5...0.....
00000020: 0001 0000 0000 0100 0000 0043 0001 0000 .....C...
...
00000070: 0000 0100 0000 0800 0000 0968 6173 6665 .....haske
00000080: 6666 3938 0000 0007 6350 5554 696d 6500 1198...CPUTime...
00000090: 0000 0462 6173 6500 0000 0847 4843 2e42 ...base...GHC.B
000000a0: 6173 6500 0000 0e47 6943 2e46 6f72 6669 ase...GHC.Forei
000000b0: 676e 5074 7200 0000 0e53 7973 7663 6a2e gnPtr....System
000000c0: 4350 5354 696d 6500 0000 0a67 6874 4330 CPUTime...getCP
000000d0: 5354 696d 6500 0000 1063 7075 5469 6865 UTime...cpuTime
000000e0: 5072 6363 6973 696f 6e Precision
```

```
type Dictionary = {count : SBH_uint32(4),
                  ids : Hstring[count]}
...
type Hi = { id : B_uint32(4) where checkId(id),
           dict : Dictionary Pointer(4), ...}
```

# Type summary

- Base types
- Tuples
- Singleton types
- Records
- Unions
- Lists/Arrays
- Value abstraction
- Type abstraction
- Dependent types
- Recursive types
- Pointers
- ???

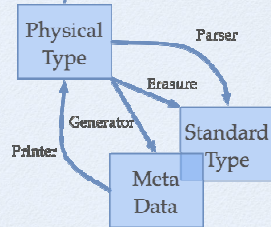
# Differences

- Data layout is not under the control of the type system.
- Physical types need some extra information: separators, terminators.
- Many physical types map to the same internal type: String('\ '), String('\'), SBH\_uint32, B\_uint32 ...
- Dependent types much more important for physical types:
- Missing value representations, value-level constraints, embedded array lengths, union tags.
- We should not assume data conforms 100% to description.

## Meta data

```
"PBAJEBG",197713,-1.948233,Quadrillion Bkn,4
"PBAJEBG",197813,-0.336838,Quadrillion Bkn,4
"PBAJEBG",197913,-1.649903,Quadrillion Bkn,4
"PBAJEBG",198013,-1.0537,Quadrillion Bkn,4
```

Described by



## Leverage

- Convert PADS description into a collection of tools:
  - Accumulators
  - Histograms
  - Clustering tool
  - Formatters
  - Translator into XML, with corresponding XML Schema.
  - XQueries using Galax's data interface
  - ...
- Long term goal:** Provide a compelling suite of tools to overcome inertia of a new language and system.

## Accumulators

- Statistical profile of "leaves" in a data source:

```
<top>.length : uint32
good: 53544 bad: 3824 pcnt-bad: 6.666
min: 35 max: 248591 avg: 4090.234
top 10 values out of 1000 distinct values:
tracked 99.552% of values
val: 3062 count: 1254 %-of-good: 2.342
val: 170 count: 1148 %-of-good: 2.144
...
SUMMING count: 9655 %-of-good: 18.032
```

Not all lengths were legal!

- Suggested by AT&T user to get "bird's eye" view of their 4000 daily feeds.
- Used at AT&T for vetting data (and for debugging PADS descriptions).

## Pretty printer

- Customizable program to reformat data:

```
207.136.97.49 - - [15/Oct/1997:18:46:51 -0700] "GET /tk/p.txt HTTP/1.0" 200 30
tj62.aol.com - - [16/Oct/1997:14:32:22 -0700] "POST /supt/ddgrp.org/confirm HTTP/1.0" 200 941
```

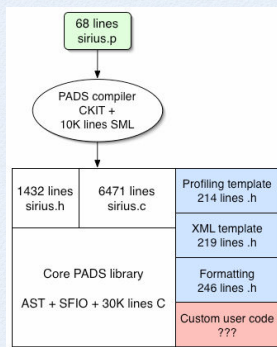
Normalize time zones  
 Normalize delimiters

Drop unnecessary values  
 Filter/repair errors

```
207.136.97.49 [-110/16/97:01:46:51]GET /tk/p.txt[1|0|200|30
tj62.aol.com[-110/16/97:21:32:22]POST /supt/ddgrp.org/confirm[1|0|200|941
```

- Users can override printing on a per type basis.
- Used at AT&T to normalize monitoring data before loading into a relational database.

## Implementation



## Why a domain specific language?

- Dramatically **shorter code** (68 lines versus ~7.9K lines).
- Description is short enough to **serve as documentation**.
- Safer:** error code inserted automatically and completely (we just have to get the compiler right...).
- We can leverage the declarative specification to produce **value-added tools**.

## Future research directions

- Design
  - How can we specify error-aware data transformations?
  - Can we infer a data transformation between two descriptions?
  - How can we express application-specific information?
  - How can we generate template programs for arbitrary data source?
- Implementation
  - How can we specialize generated libraries to incorporate application-specific information?
  - How can we optimize streaming XQueries?
- Theory
  - How do we precisely specify the semantics of PADS? (...tomorrow...)
  - What is the expressiveness of PADS vs. context free grammars?
- Engineering
  - How do we build the system to make it easy to add new base types?
  - New libraries and tools? New language bindings?

## Contributors

- Kathleen Fisher (AT&T)
- Robert Gruber (Google)
- Mary Fernandez (AT&T)
- Joel Gottlieb (AT&T)
- Yitzhak Mandelbaum (Princeton → AT&T)
- Martin Strauss (University of Michigan)
- David Walker (Princeton)
- Xuan Zheng (University of Michigan)

## Try it!

- Available for download with an open source license (CPL 1.0).
- Demo of accumulators, format program, and XML conversion.
- Send us feedback!



[www.padsproj.org](http://www.padsproj.org)