



The Next 700 Data Description Languages

Kathleen Fisher
AT&T Labs Research

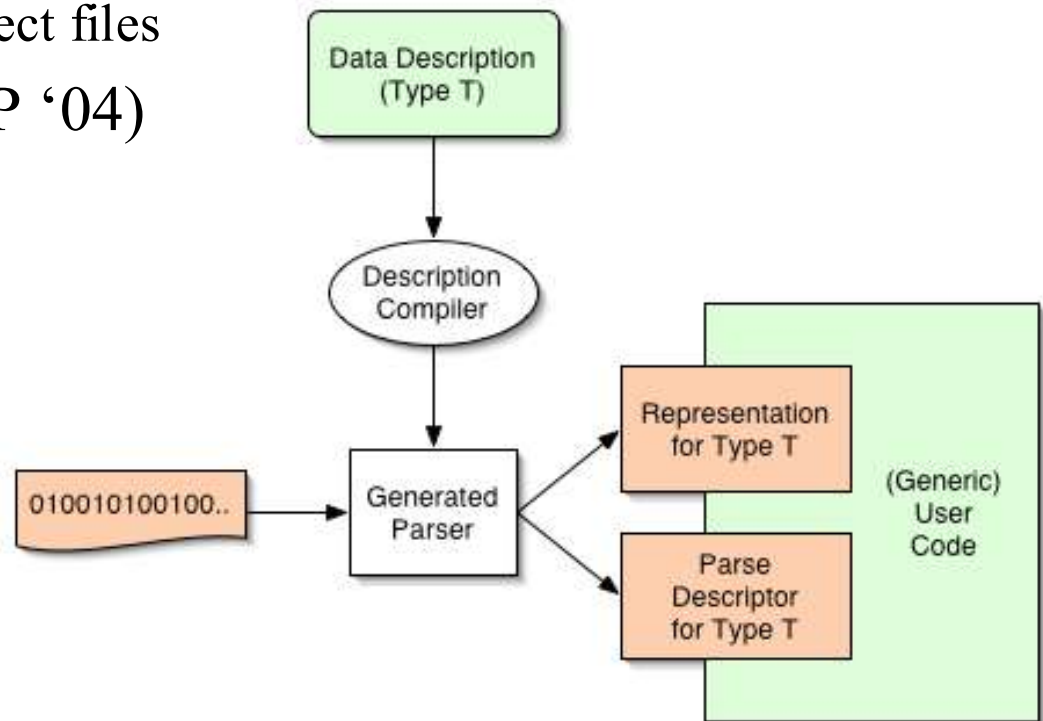
Yitzhak Mandelbaum, David Walker
Princeton

Review: Technical Challenges of Ad Hoc Data

- 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- Documentation is often out-of-date or nonexistent.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.

Many Data Description Languages

- PacketTypes (SIGCOMM '00)
 - Packet processing
- DataScript (GPCE '02)
 - Java jar files, ELF object files
- Erlang Binaries (ESOP '04)
 - Packet processing
- PADS (PLDI '05)
 - General ad hoc data



The Next 700 Programming Languages

The languages people use to communicate with computers **differ in their intended aptitudes**, towards either a particular application area, or a particular phase of computer use (high level programming, program assembly, job scheduling, etc.). They also differ **in physical appearance**, and more important, **in logical structure**. *The question arises, do the idiosyncrasies reflect basic logical properties of the situation that are being catered for? Or are they accidents of history and personal background that may be obscuring fruitful developments?* This question is clearly important if we are trying to predict or influence language evolution.

Continued...

The Next 700 Programming Languages, cont.

To answer it we must think in terms, not of languages, but **families of languages**. That is to say we must systematize their design so that a new language is a point chosen from a **well-mapped space**, rather than a laboriously devised construction.

— J. P. Landin

The Next 700 Programming Languages, 1965.

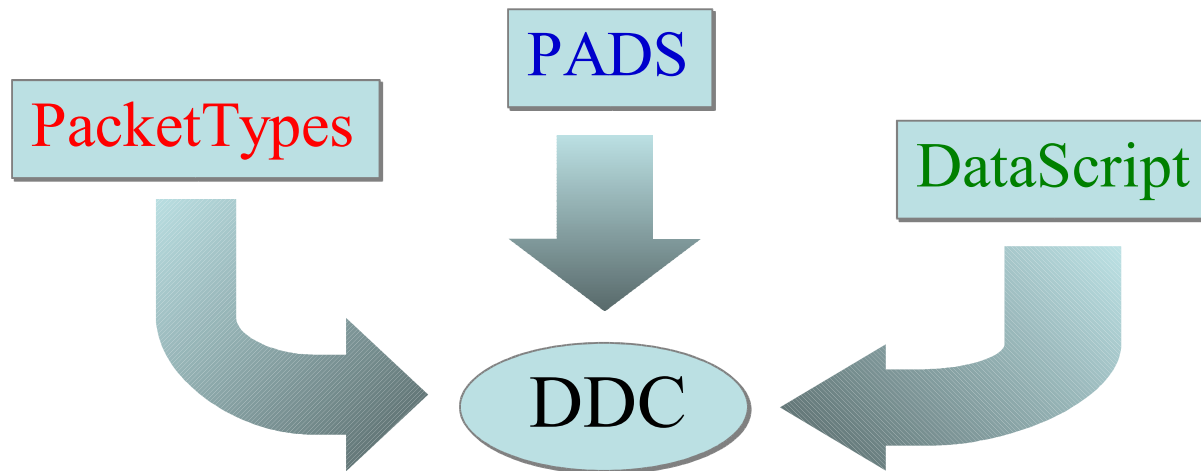
The Next 700 Data Description Languages

- What is the family of data description languages?
- How do existing languages relate to each other?
- What differences are crucial, which “accidents of history”?
- What do the existing languages mean, precisely?

To answer these questions, we introduce a semantic framework for understanding data description languages.

Contributions

- A core data description calculus (DDC)
 - Based on dependent type theory
 - Simple, orthogonal, composable types
 - Types transduce external data source to internal representation.
- Encodings of high-level DDLs in low-level DDC



Outline

- Introduction
- A Data Description “Calculus” (DDC)
- But what does DDC mean?
 - Well-kinding judgment
 - Representation, parse descriptor, and parser generation
- But what do data description languages (DDLs) mean?
 - Idealized PADS (IPADS)
 - Features from other DDLs.
- Applications of the semantics

A Data Description Calculus



Candidate DDC Primitives

- Base types parameterized by expressions ($Pstring(:|:)$)
 - Type constructor constants
- Pair of fields with cascading scope ($Pstruct$)
 - Dependent sums
- Additional constraints ($Ptypedef$, $Pwhere$, field constraints).
 - Set types
- Alternatives ($Punion$, $Popt$)
 - Sums
- Open-ended sequences ($Parray$)
 - Some kind of list?
- User-defined parameterized types
 - Abstraction and application
- “Active types”: compute, absorb, and scanning
 - Built-in functions

Base Types and Sequences

- $C(e)$: base type parameterized by expression e .
- $\Sigma x: \tau. \tau'$: dependent sum describes sequence of values.
 - Variable x gives name to first value in sequence.
 - Note syntactic sugar: $\tau * \tau'$ if x not in τ' .
- Examples:

<code>“123hello ”</code>	<code>int * string(' ') * char</code>	<code>(123, “hello”, ' ')</code>
<code>“3513”</code>	<code>$\Sigma width: int_fw(1). int_fw(width)$</code>	<code>(3, 513)</code>
<code>“:hello:”</code>	<code>$\Sigma term: char. string(term) * char$</code>	<code>(‘:’, “hello”, ‘:’)</code>

Constraints

- $\{x: \tau \mid e\}$: set types add constraints to the type τ and express relationships between elements of the data.
- Examples:

'a'	$\{c : \text{char} \mid c = \text{'a'}\}$ (<i>abbrev:</i> $S_c(\text{'a'})$)	inl 'a'
"101", "82"	$\{x : \text{int} \mid x > 100\}$	inl 101, inr 82
"43 105 67"	$\Sigma_{min:\text{int}}.S_c(\text{' '}) *$ $\Sigma_{max:\{m:\text{int} \mid min \leq m\}}.S_c(\text{' '}) *$ $\{mid:\text{int} \mid min \leq mid \ \& \ mid \leq max\}$	(43, inl ' ', inl 105, inl ' ', inl 67)

Unions and the Empty String

- $\tau + \tau'$: deterministic, exclusive or
 - try τ ; on failure, try τ' .
- **unit**: matches the empty string.
- Examples:

“54”, “n/a”	$\text{int} + S_s(\text{“n/a”})$	$\text{inl } 54, \text{inr } (\text{inl } \text{“n/a”})$
“2341”, “”	$\text{int} + \text{unit}$	$\text{inl } 2341, \text{inr } ()$

Array Features

- What features do we need to handle data sequences?
 - Elements
 - Separator between elements
 - Termination condition (“Are we done yet?”)
 - Terminator *after* sequence
- Examples:
 - “192.168.1.1”
 - “Harry|Ron|Hermione|Ginny;”

Bottom and Arrays

- $\tau \text{ seq}(\tau_s; e, \tau_t)$ specifies:
 - Element type τ
 - Separator types τ_s .
 - Termination condition e .
 - Terminator type τ_t .
- **bottom**: reads nothing, flagging an error.
- Example: IP address.

“192.168.1.1”	<code>int seq(S_c('.'); len 4, bottom)</code>	(4, [192,168,1,1])
---------------	---	--------------------

Abstraction and Application

- Can parameterize types over values: $\lambda x. \tau$
- Correspondingly, can apply types to values: τe
- Example: IP address with terminator

<i>none</i>	$\lambda term. \text{int seq}(S_c('.'); \text{len } 4, S_c(term))$	<i>none</i>
“1.2.3.4 ”	$\text{IP_addr } ' ' * S_c(' ')$	$((4, [1,2,3,4]), \text{inl } ' ')$

Absorb, Compute and Scan

- Absorb, Compute and Scan are *active* types.
 - $\text{absorb}(\tau)$: consume data from source; produce nothing.
 - $\text{compute}(e:\sigma)$: consume nothing; output result of computation e .
 - $\text{scan}(\tau)$: scan data source for type τ .
- Examples:

“ ”	$\text{absorb}(S_c())$	()
“10 12”	$\Sigma \text{width}:\text{int}.S_c() *$ $\Sigma \text{length}:\text{int}.$ $\text{area}:\text{compute}(\text{width} \times \text{length}:\text{int})$	(10,12,120)
“^%\$!&_ ”	$\text{scan}(S_c())$	(6,inl ' ')

DDC Example: Idealized Web Server Log

$S = \lambda ch. \{c : \text{char} \mid c = ch\}$

$\text{authid_t} = S('-') + \text{string}('')$

$\text{response_t} = \lambda x. \{y : \text{int16_fw}(x) \mid 100 \leq y \text{ and } y < 600\}$

$\text{entry_t} =$

$\Sigma \text{client} : \text{ip}. \quad S('')^*$

$\Sigma \text{remoteid} : \text{authid_t}. \quad S('')^*$

$\Sigma \text{response} : \text{response_t } 3.$

$\text{compute}(\text{getdomain } \text{client} = \text{"edu"} : \text{bool})$

$\text{entry_t seq}(S('\n'); \lambda x. \text{false}, \text{bottom})$

124.207.15.27 - 234
12.24.20.8 kfisher 208

A Data Description Calculus

$C(e)$	Atomic type parameterized by expression e
$\Sigma x: \tau. \tau'$	Field sequence with cascading scope
$\{x: \tau \mid e\}$	Adding constraints to existing descriptions
$\tau + \tau'$	Alternatives
$\tau \text{ seq}(\tau_s; e, \tau_d)$	Open ended sequences
$\lambda x. \tau$	Parameterizing types by expressions.
τe	Supplying values to parameterized types.
unit/bottom	Empty strings: ok/error
absorb, compute, scan	“Active types”

Semantics Overview

- Well formed DDC type: $\Gamma \vdash \tau : \kappa$
- Representation for type $\tau : [\tau]_{\text{rep}}$
- Parse descriptor for type $\tau : [\tau]_{\text{pd}}$
- Parsing function for type $\tau : [\tau]$
 - $[\tau] : \text{bits} * \text{offset} \rightarrow \text{offset} * [\tau]_{\text{rep}} * [\tau]_{\text{pd}}$

Type Kinding

Kinding ensures types are well formed.

$$\kappa ::= T \mid \sigma \rightarrow \kappa$$

$$\frac{\Gamma \vdash \tau : \sigma \rightarrow \kappa \quad \Gamma \vdash e : \sigma}{\Gamma \vdash \tau e : \kappa}$$

$$\frac{\Gamma \vdash \tau : T \quad \Gamma \vdash \tau' : T}{\Gamma \vdash \tau + \tau' : \text{type}}$$

$$\frac{\Gamma \vdash \tau : T \quad \Gamma, x : [\tau]_{\text{rep}} * [\tau]_{\text{pd}} \vdash e : \text{bool}}{\Gamma \vdash \{x : \tau \mid e\} : T}$$

Selected Representation Types

<i>DDC</i>	<i>Host Language</i>
$[C(e)]_{\text{rep}}$	$\mathbb{I}(C) + \text{none}$
$[\Sigma x: \tau. \tau']_{\text{rep}}$	$[\tau]_{\text{rep}} * [\tau']_{\text{rep}}$
$[\{x: \tau \mid e\}]_{\text{rep}}$	$[\tau]_{\text{rep}} + [\tau]_{\text{rep}}$
$[\tau + \tau']_{\text{rep}}$	$[\tau]_{\text{rep}} + [\tau']_{\text{rep}}$
$[\tau \text{ seq}(\tau_s; e, \tau_t)]_{\text{rep}}$	$\text{int} * ([\tau]_{\text{rep}} \text{ seq})$
$[\lambda x. \tau]_{\text{rep}}, [\tau e]_{\text{rep}}$	$[\tau]_{\text{rep}}$
$[\text{unit}]_{\text{rep}}$	unit

unrecoverable
error

semantic error

Note that we erase
all dependencies.

Selected Parse Descriptor Types

<i>DDC</i>	<i>Host Language</i>
$[C(e)]_{pd}$	pd_hdr
$[\Sigma x: \tau. \tau']_{pd}$	pd_hdr * $[\tau]_{pd}$ * $[\tau']_{pd}$
$[\{x: \tau \mid e\}]_{pd}$	pd_hdr * $[\tau]_{pd}$
$[\tau + \tau']_{pd}$	pd_hdr * ($[\tau]_{pd} + [\tau']_{pd}$)
$[\tau \text{ seq}(\tau_s; e, \tau_t)]_{pd}$	pd_hdr * int * int * ($[\tau]_{pd} \text{ seq}$)
$[\lambda x. \tau]_{pd}, [\tau e]_{pd}$	$[\tau]_{pd}$
$[\text{unit}]_{pd}$	pd_hdr

pd_hdr =
int * errcode * span

Parsing Semantics of Types

- Semantics expressed as parsing functions written in the polymorphic λ -calculus.
 - $[\tau] : \text{bits} * \text{offset} \rightarrow \text{offset} * [\tau]_{\text{rep}} * [\tau]_{\text{pd}}$
- Dependent sum case:

$$[\Sigma x: \tau. \tau'] =$$
$$\lambda(B, \omega).$$
$$\text{let } (\omega_1, r_1, p_1) = [\tau] (B, \omega) \text{ in}$$
$$\text{let } x = (r_1, p_1) \text{ in}$$
$$\text{let } (\omega_2, r_2, p_2) = [\tau'] (B, \omega_1) \text{ in}$$
$$(\omega_2, R_{\Sigma}(r_1, r_2), P_{\Sigma}(p_1, p_2))$$

Properties of the Calculus

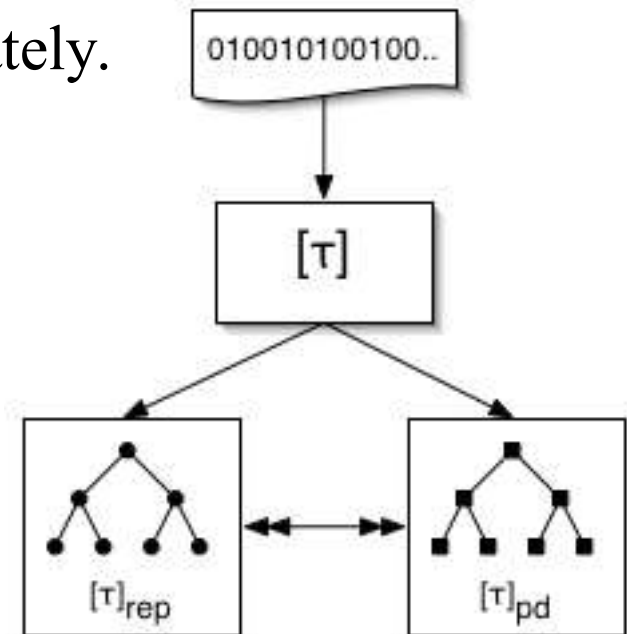
- **Theorem:** If $\Gamma \vdash \tau : \kappa$ then

$$\Gamma \vdash [\tau] : \text{bits} * \text{offset} \rightarrow \text{offset} * [\tau]_{\text{rep}} * [\tau]_{\text{pd}}$$

“Well-formed type τ yields a parser that returns values with types corresponding to τ .”

- **Theorem:** Parsers report errors accurately.

- Errors in parse descriptor correspond to errors in representation.
- Parsers check all semantic constraints.



Making Use of the Calculus

IPADS

IPADS

$t ::= \mathbf{C}(e) \mid \mathbf{Pfun}(x:s) = t \mid t e$
| $\mathbf{Pstruct}\{\text{fields}\} \mid \mathbf{Punion}\{\text{fields}\}$
| $\mathbf{Pswitch} e \text{ of } \{\text{alts } t_{\text{def}i}\} \mid \mathbf{Popt} t$
| $t \mathbf{Pwhere} x.e \mid \mathbf{Palt}\{\text{fields}\}$
| $t \mathbf{Parray} [t, t] \mid \mathbf{Pcompute} e \mid \mathbf{Plit} c$

$\text{fields} ::= \mid \text{fields } x : t;$
 $\text{alts} ::= \mid \text{alts } e \Rightarrow t;$

DDC

$t \Rightarrow \tau$

IPADS Example

```
authid_t = Punion { unauth : "-"; id : Pstring(" ");};
```

```
response_t = Pfun(x:int)
```

```
    Puint16_FW(x) Pwhere y.100 <= y and y < 600;
```

```
entry_t = Pstruct {
```

```
    client : Pip;           " ";
```

```
    remoteid : authid_t;    " ";
```

```
    response : response_t 3;
```

```
    academic : Pcompute(getdomain client = "edu" : bool);
```

```
};
```

```
entry_t Parray(Peor, Peof)
```

```
124.207.15.27 - 234  
12.24.20.8 kfisher 208
```

Example: Popt and Plit

unit
 $\tau_1 + \tau_2$

$t \Rightarrow \tau$

Popt $t \Rightarrow \tau + \text{unit}$

$C(e)$
 $\{x:\tau \mid e\}$
absorb(τ)
scan(τ)

$c : \text{char}$

Plit $c \Rightarrow \text{scan}(\text{absorb}(\{x:\text{char} \mid x = c\}))$

Example: Pswitch

$$\begin{array}{l} \tau + \tau' \\ \lambda x. \tau \\ \{x:\tau | e\} \end{array}$$

$$t_i \Rightarrow \tau_i \quad (i = 1 \dots n)$$

$$t_{\text{def}} \Rightarrow \tau_{\text{def}}$$

Pswitch e of $\{e_1 \Rightarrow t_1; e_2 \Rightarrow t_2; \dots t_{\text{def}}\} \Rightarrow$

$$(\lambda c. \{x:\tau_1 | c = e_1\} + \{x:\tau_2 | c = e_2\} + \dots + \tau_{\text{def}}) e$$

Example: Pswitch

$\tau + \tau'$
$\lambda x. \tau$
$\{x:\tau e\}$

$$t_i \Rightarrow \tau_i \quad (i = 1 \dots n)$$

$$t_{\text{def}} \Rightarrow \tau_{\text{def}}$$

Pswitch e of $\{e_1 \Rightarrow t_1; e_2 \Rightarrow t_2; \dots t_{\text{def}}\} \Rightarrow$

$$(\lambda c. \{x:\tau_1 | c = e_1\} + \{x:\tau_2 | c = e_2\} + \dots + \tau_{\text{def}}) e$$

But this encoding isn't exactly right, as it parses the data as each branch until it reaches the matching tag.

Encoding Conditionals

if e then t_1 else t_2

$$t_1 \Rightarrow \tau_1 \qquad t_2 \Rightarrow \tau_2$$

$$\text{if } e \text{ then } t_1 \text{ else } t_2 \Rightarrow (\{x:\text{unit} \mid !e\} + \tau_1) * (\{x:\text{unit} \mid e\} + \tau_2)$$

Pswitch Revisted

- Encode **Pswitch** as a sequence of conditionals

<pre>Pswitch e { e₁ => x₁ : t₁ ... e_n => x_n : t_n t_{def} }</pre>	=	<pre>(Pfun (x : int) = if x = e₁ then t₁ else ... if x = e_n then t_n else t_{def}) e</pre>
---	---	--

Other Features

- PacketTypes: arrays, where clauses, structures, overlays, and alternation.
- DataScript: set types (enumerations and bitmask sets), arrays, constraints, value-parameterized types, and (monotonically increasing labels).

Other Uses of the Semantics

- Bug hunting!
 - Non-termination of array parsing if no progress made.
 - Inconsistent parse descriptor construction.
- Principled extensions
 - Adding recursion (done)
 - Adding polymorphism?
- Distinguishing the essential from the accidental
 - Highlights places where PADS sacrifices safety.
 - **Pomit** and **Pcompute** : much more useful than originally thought
 - **Punion** : what if correct branch has an error?

Future work

- What are the set of languages recognized by the DDC?
- How does the expressive power of the DDC relate to CFGs and regular expressions?
- Add polymorphism to DDC and PADS.

Summary

- Data description languages are well-suited to describing ad hoc data.
- No one DDL will ever be right - different domains and applications will demand different languages with differing levels of expressiveness and abstraction.
- Our work defines the first semantics for data description languages.
- For more information, visit www.padsproj.org.