



The Next 700 Data Description Languages

Kathleen Fisher
AT&T Labs Research
Yitzhak Mandelbaum, David Walker
Princeton

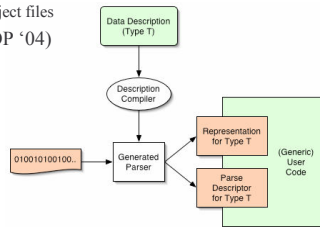
Review: Technical Challenges of Ad Hoc Data

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- Documentation is often **out-of-date** or **nonexistent**.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

Summer School 2007

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



Summer School 2007

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...

Summer School 2007

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.

Summer School 2007

The Next 700 Data Description Languages

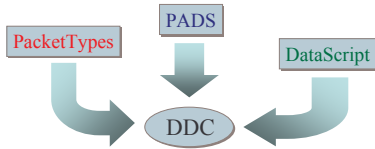
- What is the family of data description languages?
- How do existing languages relate to each other?
- What differences are crucial, which “accidents of history”?
- What do the existing languages mean, precisely?

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

Summer School 2007

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



Summer School 2007

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

Summer School 2007

A Data Description Calculus



Summer School 2007

Candidate DDC Primitives

- Base types parameterized by expressions (Pstring(' : '))
 - Type constructor constants
- Pair of fields with cascading scope (Pstruct)
 - Dependent products
- 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

Summer School 2007

Base Types and Sequences

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

"123hello"	$\text{int} * \text{string}(' ') * \text{char}$	(123, "hello", ' ')
"3513"	$\Sigma \text{width}: \text{int}. \text{fw}(1). \text{int_fw}(\text{width})$	(3, 513)
"hello:"	$\Sigma \text{term}: \text{char}. \text{string}(\text{term}) * \text{char}$	(';', "hello", ':')

Summer School 2007

Constraints

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

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

Summer School 2007

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"	int + $S_c("n/a")$	inl 54, inr "n/a"
"2341", ""	int + unit	inl 2341, inr ()

Summer School 2007

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;"

Summer School 2007

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"	int seq($S_c('.')$; len 4, bottom)	[192,168,1,1]
---------------	--------------------------------------	---------------

Summer School 2007

Abstraction and Application

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

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

Summer School 2007

Absorb, Compute and Scan

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

" "	absorb($S_c(' ')$)	()
"10 12"	$\Sigma width:\text{int}. S_c(' ') *$ $\Sigma length:\text{int}.$ $area:\text{compute}(width \times length:\text{int})$	(10,12,120)
"^0%\$1&_ "	scan($S_c(' ')$)	(6,inl ' ')

Summer School 2007

DDC Example: Idealized Web Server Log

$S = \lambda st. \{s : \text{string} \mid s = st\}$
 $\text{authid}_t = S(" ") + \text{string}('')$
 $\text{response}_t = \lambda x. \{y : \text{int16_FW}(x) \mid 100 \leq y \text{ and } y < 600\}$
 $\text{entry}_t =$

- $\Sigma client : \text{ip}.$ $S(" ") *$
- $\Sigma remoteid : \text{authid}_t.$ $S(" ") *$
- $\Sigma response : \text{response}_t 3.$
- compute(getdomain $client = "edu"$: bool)

 $\text{entry}_t \text{ seq}(S("n")). \lambda x. \text{false, bottom}$

124.207.15.27 - 234
12.24.20.8 kfisher 208

Summer School 2007

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"

Summer School 2007

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 $\tau : [\tau]$
 - $[\tau] : \text{bits} * \text{offset} \rightarrow \text{offset} * [\tau]_{\text{rep}} * [\tau]_{\text{PD}}$

Summer School 2007

Type Kinding

- Kinding ensures types are well formed.

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

$$\frac{\Gamma \vdash \tau : \text{type} \quad \Gamma \vdash \tau' : \text{type}}{\Gamma \vdash \tau + \tau' : \text{type}}$$

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

Summer School 2007

Selected Representation Types

DDC	Host Language
$[C(e)]_{\text{rep}}$	$\mathbb{I}(C) + \text{noval}$ ← unrecoverable error
$[\Sigma x: \tau. \tau']_{\text{rep}}$	$[\tau]_{\text{rep}} * [\tau']_{\text{rep}}$
$[\{x: \tau \mid e\}]_{\text{rep}}$	$[\tau]_{\text{rep}} + ([\tau]_{\text{rep}} \text{ error})$ ← semantic error
$[\tau + \tau']_{\text{rep}}$	$[\tau]_{\text{rep}} + [\tau']_{\text{rep}} + \text{noval}$
$[\tau \text{ seq}(\tau_s; e, \tau_d)]_{\text{rep}}$	$\text{int} * ([\tau]_{\text{rep}} \text{ seq})$
$[\lambda x. \tau]_{\text{rep}}, [\tau e]_{\text{rep}}$	$[\tau]_{\text{rep}}$
$[\text{unit}]_{\text{rep}}$	unit

Note that we erase all dependencies.

Summer School 2007

Selected Parse Descriptor Types

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

pd_hdr =
int * errcode * span

Summer School 2007

Parsing Semantics of Types

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

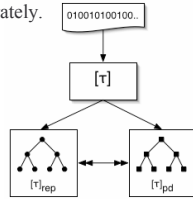
- Product case:

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

Summer School 2007

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.



Summer School 2007

Making Use of the Calculus

IPADS

$$t ::= C(e) \mid Pfun(x:s) = t \mid t e$$

$$\mid Pstruct\{fields\} \mid Punion\{fields\}$$

$$\mid Pswitch e \text{ of } \{alts\ t_{def}\} \mid Popt t$$

$$\mid t Pwhere\ x.e \mid Palt\{fields\}$$

$$\mid t Parray [t, t] \mid Pcompute e \mid Pplit c$$

fields ::= | fields x : t;
alts ::= | alts e => t;



Summer School 2007

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

Summer School 2007

Example: Popt and Pplit

unit
 $\tau_1 + \tau_2$

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

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

$$\frac{c : char}{Pplit\ c \Rightarrow scan(absorb(\{x:char \mid x = c\}))}$$

Summer School 2007

Example: Pswitch

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

$$\frac{t_i \Rightarrow \tau_i \ (i = 1 \dots n) \quad t_{def} \Rightarrow \tau_{def}}{Pswitch\ e\ \text{of}\ \{e_1 \Rightarrow t_1; e_2 \Rightarrow t_2; \dots t_{def}\} \Rightarrow (\lambda c. \{x:\tau_1 \mid c = e_1\} + \{x:\tau_2 \mid c = e_2\} + \dots + \tau_{def}) e}$$

Summer School 2007

Example: Pswitch

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

$$\frac{t_i \Rightarrow \tau_i \ (i = 1 \dots n) \quad t_{def} \Rightarrow \tau_{def}}{Pswitch\ e\ \text{of}\ \{e_1 \Rightarrow t_1; e_2 \Rightarrow t_2; \dots t_{def}\} \Rightarrow (\lambda c. \{x:\tau_1 \mid c = e_1\} + \{x:\tau_2 \mid c = e_2\} + \dots + \tau_{def}) e}$$

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

Summer School 2007

Encoding Conditionals

if e then t_1 else t_2

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

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

Summer School 2007

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>
---	---	--

Summer School 2007

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).

Summer School 2007

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 (done in PADS/ML)
- Distinguishing the essential from the accidental
 - Highlights places where PADS/C sacrifices safety.
 - **Pomit** and **Pcompute** : much more useful than originally thought
 - **Punion** : what if correct branch has an error?

Summer School 2007

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.

Summer School 2007