

The Hadoop Pig Syntax Card



Types

<code>int</code>	32-bit signed integer	<code>42</code>
<code>long</code>	64-bit signed integer	<code>42L</code>
<code>float</code>	32-bit signed floating point	<code>3.14F 0.0314e2F</code>
<code>double</code>	64-bit signed floating point	<code>2.718 271.8e-2</code>
<code>chararray</code>	UTF-8 character string	<code>'Hadoop'</code>
<code>bytearray</code>	array of bytes	<code>N/A</code>
<code>tuple</code>	ordered set of fields	<code>(field [, field ...])</code>
<code>bag</code>	unordered collection of tuples	<code>{tuple [, tuple ...]}</code>
<code>map</code>	set of key value pairs	<code>[key#val [, key#val ...]]</code>

Schemas

Simple Types	<code>(alias[:type] [, alias[:type] ...])</code>
Tuple	<code>(alias[:tuple] (alias[:type] [, alias[:type] ...]))</code>
Bag	<code>(alias[:bag] { tuple schema })</code>
Map	<code>(alias[:map] [])</code>

Dereference

Tuple	<code>T.field_name</code> <code>T.\$0</code> <code>T.(\$1, field_name)</code>
Bag	<code>B.field_name</code> <code>B.\$0</code> <code>B.(\$1, field_name)</code>
Map	<code>M#'key'</code>

UDFs

```
REGISTER path;

DEFINE alias {function|['command' [input] [output] [ship] [cache]]};

input  INPUT {stdin|'path'} [USING serializer]
        [, {stdin|'path'} [USING serializer] ...]

output OUTPUT {stdin|stderr|'path'} [USING serializer]
        [, {stdin|stderr|'path'} [USING serializer] ...]

ship   SHIP ( 'path' [, 'path' ... ] )

cache  CACHE ( 'dfs_path#dfs_file' [, 'dfs_path#dfs_file' ... ] )
```

Eval	Math	String	Bag/Tuple
AVG	ABS	INDEXOF	TOBAG
CONCAT	ACOS	LAST_INDEX_OF	TOP
COUNT	ASIN	LCFIRST	TOTUPLE
COUNT_STAR	ATAN	LOWER	
DIFF	CBRT	REGEX_EXTRACT	
MAX	CEIL	REGEX_EXTRACT_ALL	
MIN	COS	REPLACE	
SIZE	COSH	STRSPLIT	
SUM	EXP	SUBSTRING	
TOKENIZE	FLOOR	TRIM	
		UCFIRST	
		UPPER	

Misc

```
-x {local,mapreduce}    Choose execution mode
-param param=value      Assign value to param for use in the script
-param_file file        Load parameters from file
$param                  Dereference parameter param
%declare param value    Assign value to param in the script
%default param value     Set default value for param
```

```
fs FSShell_subcommand parameters
sh shell_subcommand parameters
{exec,run} [-param name=value] [-param_file file] script
kill jobid
set key 'value'
help
quit
```



The Hadoop Pig Syntax Card



Input/Output

```
alias = LOAD 'data' [USING function] [AS schema];  
  
STORE alias INTO 'directory' [USING function];  
  
alias = ORDER alias BY { * [ASC|DESC]  
                        | field_alias [ASC|DESC]  
                        | [, field_alias [ASC|DESC] ...]  
                        }  
  [PARALLEL n];
```

Diagnostic

```
DESCRIBE alias;  
  
DUMP alias;  
  
EXPLAIN alias;
```

Selection

```
alias_l = LIMIT alias_O n;  
  
alias_l = SAMPLE alias_O size; (size in [0.0, 1.0])  
  
alias_l = DISTINCT alias_O;  
  
alias_l = FILTER alias_O BY expression;  
  
SPLIT alias_O INTO alias_l IF expression_l  
  [, alias_N IF expression_N ...];
```

Joining and Grouping

```
alias = UNION [ONSCHEMA] alias_O, alias_l [, alias_N ...];  
  
alias = CROSS alias_O, alias_l [, alias_N ...] [PARTITION BY part] [PARALLEL n];  
  
alias = JOIN alias_O BY {expression|('expression [, expression ...]')}  
  [, alias_l BY {expression|('expression [, expression ...]')} ...]  
  [USING 'replicated' | 'skewed' | 'merge']  
  [PARTITION BY partitioner] [PARALLEL n];  
  
alias = JOIN left_alias BY left_alias_column [LEFT|RIGHT|FULL] [OUTER],  
  right_alias BY right_alias_column  
  [USING 'replicated' | 'skewed' | 'merge']  
  [PARTITION BY partitioner] [PARALLEL n];  
  
alias = GROUP alias { ALL | BY expression } [, alias ALL | BY expression ...]  
  [USING 'collected' | 'merge']  
  [PARTITION BY partitioner] [PARALLEL n];
```

Transformation

```
alias_l = FOREACH alias GENERATE expression [AS schema]  
  [, expression [AS schema] ... ];  
  
alias_l = FOREACH alias {  
  alias = nested_op; [ alias = nested_op; ... ]  
  GENERATE expression [AS schema][, expression [AS schema] ... ]  
};  
  
alias_l = MAPREDUCE 'mr.jar'  
  STORE alias_2 INTO 'inputLocation' USING storeFunc  
  LOAD 'outputLocation' USING loadFunc AS schema [ `params, ... ` ];  
  
alias = STREAM alias [, alias ...]  
  THROUGH { `command` | cmd_alias } [AS schema];
```

