

Assignment 1

Adding a Set Data Type to PostgreSQL

Last updated: **Tuesday 9th March 11:05pm**

Most recent changes are shown in **red**;
older changes are shown in **brown**.

This is a working draft.
Small updates and clarifications will be added as needed.
No major changes are anticipated.

[jump to ...](#) [summary](#) [introduction](#) [setup](#) [intsets](#) ([values](#) [operations](#) [representing](#)) [changelog](#) [submission](#)

Summary

Deadline	Friday 19 March, 9:00pm
Marks	This assignment contributes 15 marks toward your total mark for this course.
Late Penalty	<i>0.089 marks</i> off the ceiling mark for each hour late. After 7 <i>days</i> the assignment is worth <i>0 marks</i> .
Submission	From WebCMS: COMP9315 website , -> <i>Assignments</i> , -> <i>Ass1 Specification</i> , -> <i>Make Submission</i> , -> upload <code>intset.c</code> , <code>intset.source</code> , [<i>other .c files</i> > <i>Makefile</i>] or From CSE: give cs9315 ass1 intset.c intset.source [<i>other .c files</i>> <i>Makefile</i>]
Pre-requisites	Before starting this assignment, it would be useful to complete Prac Exercise P04.

This assignment aims to give you:

- An understanding of how data is treated inside a DBMS
- Practice in adding a new base type to PostgreSQL

The goal is to implement a new data type for PostgreSQL, complete with input/output functions and a range of operations.

Make sure that you read this assignment specification *carefully* and *completely* before starting work on the assignment. Questions which indicate that you haven't done this will simply get the response "Please read the spec".

We will be using the following names in the discussion below:

PG_CODE	The directory where your PostgreSQL source code is located (typically <code>/srvr/\$USER/postgresql-12.5</code>)
PG_HOME	The directory where you have installed the PostgreSQL binaries (typically <code>/srvr/\$USER/pgsql</code>)
PG_DATA	The directory where you have placed PostgreSQL's data (typically <code>/srvr/\$USER/pgsql/data</code>)

PG_LOG the file to where you send PostgreSQL's log output
(typically `/srvr/$USER/pgsql/log`)

Introduction

PostgreSQL has an extensibility model which, among other things, provides a well-defined process for adding new data types into a PostgreSQL server.

This capability has led to the development by PostgreSQL users of a number of types (such as polygons) which have become part of the standard distribution.

It also means that PostgreSQL is the database of choice in research projects which aim to push the boundaries of what kind of data a DBMS can manage.

In this assignment, we will be adding a new data type for dealing with **sets of integers**.

You may implement the functions for the data type in any way you like, *provided that* they satisfy the semantics given below.
(in the `intSets` section)

Note that arrays in PostgreSQL have some properties and operations that make them look a little bit like sets. However, they are *not* sets and have quite different semantics to the data type that we are asking you to implement.

The process for adding new base data types in PostgreSQL is described in the following sections of the PostgreSQL documentation:

- [37.10 C-Language Functions](#)
- [37.13 User-defined Types](#)
- [37.14 User-defined Operators](#)
- [SQL: CREATE TYPE](#)
- [SQL: CREATE OPERATOR](#)
- [SQL: CREATE OPERATOR CLASS](#)

Section 37.13 uses an example of a complex number type, which you can use as a starting point for defining your `intSet` data type (see below). Note that the complex type is a starting point *only*, to give an idea on how new data types are added. Don't be fooled into thinking that this assignment just requires you to change the name `complex` to `intSet`; the `intSet` type is more complex (no pun intended) than the complex number type.

There are other examples of new data types under the `tutorial` and `contrib` directories. These may or may not give you some useful ideas on how to implement the `intSet` data type:

- An auto-encrypted password datatype
`PG_CODE/contrib/chkpass/`
- A case-insensitive character string datatype
`PG_CODE/contrib/citext/`
- A confidence-interval datatype
`PG_CODE/contrib/seg/`

Setting Up

You ought to start this assignment with a fresh copy of PostgreSQL, without any changes that you might have made for the Prac exercises (unless these changes are trivial).

Note that you only need to configure, compile, and install your PostgreSQL server once for this assignment.

All subsequent compilation takes place in the `src/tutorial` directory, and only requires modification of the files there.

Once you have re-installed your PostgreSQL server, you should run the following commands:

```
$ cd PG_CODE/src/tutorial
$ cp complex.c intset.c
$ cp complex.source intset.source
```

Once you've made the `intset` files, you should also edit the Makefile in this directory, and add the **green** text to the following lines:

```
MODULES = complex funcs intset
DATA_built = advanced.sql basics.sql complex.sql funcs.sql syscat.sql intset.sql
```

The rest of the work for this assignment involves editing the `intset.c` and `intset.source` files.

In order for the Makefile to work properly, you must use the identifier `_OBJWD_` in the `intset.source` file to refer to the directory holding the compiled library.

You should never directly modify the `intset.sql` file produced by the Makefile.

If you want to use other `*.c` files along with `intset.c`, then you can do so, but you will need to make further changes to the Makefile to ensure that they are compiled and linked correctly into the library.

Note that your submitted versions of `intset.c` and `intset.source` should not contain any references to the `complex` type (because that's not what you're implementing).

Make sure that the comments in the program describes the code that *you* wrote.

Also, *do not* put testing queries in your `intset.source`; all it should do is create the new data type. Put any testing you want to do in a separate `*.sql`, which you don't need to submit. And *do not* drop the `intSet` type at the end of `intset.source`. If you do, your data type will vanish before we have a chance to test it.

The intSet Data Type

We aim to define a new base type `intSet`, to store the notion of sets of integer values.

We also aim to define a useful collection of operations on the `intSet` type.

How you represent `intSet` values, and implement functions to manipulate them, is up to you.

However, they must satisfy the requirements below

Once implemented correctly, you should be able to use your PostgreSQL server to build the following kind of SQL applications:

```
create table Features (  
    id integer primary key,  
    name text  
);  
  
create table DBSystems (  
    name text primary key,  
    features intSet  
);  
  
insert into Features (id, name) values  
    (1, 'well designed'),  
    (2, 'efficient'),  
    (3, 'flexible'),  
    (4, 'robust');  
  
insert into DBSystems (name, features) values  
    ('MySQL', '{}'),  
    ('MongoDB', '{}'),  
    ('Oracle', '{2,4}'),  
    ('PostgreSQL', '{1,2,3,4}');
```

intSet values

In mathematics, we represent a set as a curly-bracketed, comma-separated collection of values: $\{1, 2, 3, 4, 5\}$. Such a set contains only distinct values, and no particular ordering can be imposed.

Our intSet values can be represented similarly.

We can have a comma-separated list of **non-negative** integers, surrounded by a set of curly braces, which is presented to and by PostgreSQL as a string.

For example:

'{ 1, 2, 3, 4, 5 }'.

Whitespace should not matter, so '{1,2,3}' and '{ 1, 2, 3 }' are equivalent.

Similarly, a set contains distinct values, so '{1,1,1,1,1}' is equivalent to '{1}'.

And ordering is irrelevant, so '{1,2,3}' is equivalent to '{3,2,1}'.

The integer values in the set are assumed to consist of a sequence of digits. There are no + or - signs.

There can be leading zeroes, but they should effectively be ignored, e.g. 0001 should be treated the same as 1.

You **may not assume a fixed limit** to the size of the set.

It may contain zero or more elements, bounded by the database's capacity to store values.

You **may assume** that each interger value will be less than INT_MAX.

ie. each element in the set will be less than $2^{31} - 1$.

Valid intSets

```
'{ }'
'{2,3,1}'
'{6,6,6,6,6,6}'
'{10, 9, 8, 7, 6,5,4,3,2,1}'
'{1, 999, 13, 666, 5}'
'{ 1 , 3 , 5 , 7,9 }'
'{1, 01, 001, 0001}' (same as '{1}')
'{1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20}'
'{1,2,3,4,5,6,7,8,9,10,(and then all numbers to),9999,10000}'
' {1,5,7,9}'
'{2,4,6, 8 }'
'{0}'
```

Invalid intSets

```
'{a,b,c}'
'{ a, b, c }'
'{1, 2.0, 3}'
'{1, {2,3}, 4}'
'{1, 2, 3, 4, five}'
'{ 1 2 3 4 }'
'1 2 3 4'
'1 2 3}'
'{ -1 }'
'{1,2,3}'
'1,2,3,4,5'
'{{1,2,3,5}}'
'{7,17,,27,37}'
'{1,2,3,5,8,}'
```

Operations on intSets

You must implement **all** of the following operations for the intSet type.
(assuming A , B , and S are intSets, and i is an integer):

- i ? S** intSet S contains the integer i ;
That is, $i \in S$.
- # S** Give the *cardinality*, or number of distinct elements in, intSet S ;
That is, $|S|$.
- A >@ B** Does intSet A contain all the values in intSet B ?
ie, for every element of B , is it an element of A ?
That is, **the improper superset** ($A \supseteq B$)
- A @< B** Does intSet A contain only values in intSet B ?
ie, for every element of A , is it an element of B ?
That is, **the improper subset** ($A \subseteq B$)
- A = B** intSets A and B are equal;
That is, intSet A contains all the values of intSet B and intSet B contains all the values of intSet A ;
That is, every element in A can be found in B , and vice versa.
- A <> B** intSets A and B are not equal;
That is, intSet A doesn't contain all the values of intSet B or intSet B doesn't contain all the values of intSet A ;
That is, **some element in A cannot be found in B , or vice versa.**
- A && B** Takes the *set intersection*, and produces an intSet containing the elements common to A and B ;
That is, $A \cap B$.

- A || B** Takes the *set union*, and produces an intSet containing all the elements of A and B ; That is, $A \cup B$.
- A !! B** Takes the *set disjunction*, and produces an intSet containing elements that are in A and not in B , or that are in B and not in A .
- A - B** Takes the *set difference*, and produces an intSet containing elements that are in A and not in B . Note that this is *not* the same as $A !! B$.

Below are examples of how you might use intSets, to illustrate the semantics.

You can use these as an initial test set; we will supply a more comprehensive test suite later.

```
db=# create table mySets (id integer primary key, iset intSet);
CREATE TABLE
db=# insert into mySets values (1, '{1,2,3}');
INSERT 0 1
db=# insert into mySets values (2, '{1,3,1,3,1}');
INSERT 0 1
db=# insert into mySets values (3, '{3,4,5}');
INSERT 0 1
db=# insert into mySets values (4, '{4,5}');
INSERT 0 1
db=# select * from mySets order by id;
 id | iset
----+-----
  1 | {1,2,3}
  2 | {1,3}
  3 | {3,4,5}
  4 | {4,5}
(4 rows)
-- get all pairs of tuples where the second iset is a subset of first iset
db=# select a.*, b.* from mySets a, mySets b
db=# where (b.iset @< a.iset) and a.id != b.id;
 id | iset | id | iset
----+-----+----+-----
  1 | {1,2,3} | 2 | {1,3}
  3 | {3,4,5} | 4 | {4,5}
(2 rows)
-- insert extra values into the iset in tuple #4 via union
db=# update mySets set iset = iset || '{5,6,7,8}' where id = 4;
UPDATE 1
db=# select * from mySets where id=4;
 id | iset
----+-----
  4 | {4,5,6,7,8}
```

```

(1 row)
-- tuple #4 is no longer a subset of tuple #3
db=# select a.*, b.* from mySets a, mySets b
db=# where (b.iset @< a.iset) and a.id != b.id;
 id | iset | id | iset
-----+-----+-----+-----
  1 | {1,2,3} | 2 | {1,3}
(1 row)
-- get the cardinality (size) of each intSet
db=# select id, iset, (#iset) as card from mySets order by id;
 id | iset | card
-----+-----+-----
  1 | {1,2,3} | 3
  2 | {1,3} | 2
  3 | {3,4,5} | 3
  4 | {4,5,6,7,8} | 5
(4 rows)
-- form the intersection of each pair of sets
db=# select a.iset, b.iset, a.iset && b.iset
db=# from mySets a, mySets b where a.id < b.id;
 iset | iset | ?column?
-----+-----+-----
 {1,2,3} | {1,3} | {1,3}
 {1,2,3} | {3,4,5} | {3}
 {1,2,3} | {4,5,6,7,8} | {}
 {1,3} | {3,4,5} | {3}
 {1,3} | {4,5,6,7,8} | {}
 {3,4,5} | {4,5,6,7,8} | {4,5}
(6 rows)
db=# delete from mySets where iset @< '{1,2,3,4,5,6}';
DELETE 3
db=# select * from mySets;
 id | iset
-----+-----
  4 | {4,5,6,7,8}
(1 row)
-- etc. etc. etc.

```

You should think of some more tests of your own.

In particular, make sure that you check that your code works with large `intSets` (e.g. cardinality ≥ 1000).

If you come up with any tests that you think are particularly clever, feel free to post them in the comments section below.

Representing intSets

The first thing you need to do is to decide on an internal representation for your `intSet` data type. You should do this *after* you have understood the description of the operators above. Since what they require may affect how you decide on the representation of your `intSet` values.

Note that because of the requirement that an `intSet` can be arbitrarily large (see above), you **cannot** have a representation that uses a fixed-size object to hold values of type `intSet`.

When you read strings representing `intSet` values, they are converted into your internal form, stored in the database in this form, and operations on `intSet` values are carried out using this data structure. When you display `intSet` values, you should show them in a canonical form, regardless of how they were entered or how they are stored. The canonical form for output (at least) should include no spaces, and should have elements in ascending order.

The first functions you need to write are ones to read and display values of type `intSet`. You should write analogues of the functions `complex_in()` and `complex_out()` that are defined in the file `complex.c`. Suitable names for these functions would be e.g. `intset_in()` and `intset_out()`. Make sure that you use the V1 style function interface (as is done in `complex.c`).

Note that the two input/output functions should be complementary, meaning that any string displayed by the output function must be able to be read using the input function. There is no requirement for you to retain the precise string that was used for input (e.g. you could store the `intSet` value internally in canonical form).

Note that you are *not* required to define binary input/output functions called `receive_function` and `send_function` in the PostgreSQL documentation, and called `complex_send()` and `complex_recv()` in the `complex.c` file.

Hint: test out as many of your C functions as you can *outside* PostgreSQL (e.g., write a simple test driver) Before you try to install them in PostgreSQL. This will make debugging much easier.

You should ensure that your definitions *capture the full semantics of the operators* (e.g. specify commutativity if the operator is commutative).

ChangeLog

- **v1.0** (2021-02-26 15:00:00+10:00)
 - released Assignment 1
- **v1.1** (2021-02-27 11:00:00+10:00)
 - Modify the "Operations on intSets" section
 - Change the symbol for "contains" from '<@' to '?'
 - Change the symbol for "cardinality" from '@' to '#'
 - Change the symbol for "subset" from '@>' to '@<'
 - Add the "superset" operation, using the '>@' symbol
 - Add the "inequality" operation, using the '<>' symbol

- **v1.2** (2021-02-26 16:00:00+10:00)
 - Corrected typo in examples of Valid IntSets
 - Add additional examples of (In)Valid IntSets
 - Add an upper bound for the size of each element in an IntSet
- **v1.3** (2021-02-26 18:00:00+10:00)
 - Correct superset and subset symbols used in "Operations on intSets"
- **v1.4** (2021-03-01 18:00:00+10:00)
 - Moved deadline forward to 9pm to be consistent with CSE policy

Submission

You need to submit two files:

`intset.c` - containing the C functions that implement the internals of the `intSet` data type.

`intset.source` - containing the template SQL commands to install the `intSet` data type into a PostgreSQL server.

Do *not* submit the `intset.sql` file, since it contains absolute file names which are not useful in our test environment.

If your system requires other `*.c` files, you should submit them, along with the modified `Makefile` from the `src/tutorial` directory.

Do not include:

- `create table ...`
- `insert into ...`
- `select ...`
- `drop type ...`

Or any other statements not directly needed for creating the `intSet` data type in `intset.source`.

Have fun, *jas* and *dylan*.