



LOFAR synthesis data handling: TaQL

Dwingeloo, 25-Aug-2015

Why TaQL?



One often wants to:

- Look at values in a MeasurementSet or other Casacore tables
- Check if data are correct
- Make small modifications
- Do some calculations
- Convert values to another reference frame or to readable text

Can be done in a Python script using pyrap.tables

- Can take some effort

Can often be done more easily using TaQL (Table Query Language)

- Very similar to SQL
- Easy to use for simple problems
 - (although some people think it is very complex)
- Can handle complex problems, but requires some more knowledge (and thought) you have to know the function names, etc.

Fully described in http://casacore.github.io/casacore-notes/199.html

MeasurementSet



A structured set of tables			
 Main table contains the visibility data, flags, UVW, weights, … 			
A collection of rows/columns	A collection of rows/columns where each row contains the data of a single baseline and timeslot (and band/field):		
- DATA	2D Complex array (nfreq*4); visibilities (XX,XY,YX,YY)		
- FLAG	2D Bool array; flag per visibility (True=bad)		
- UVW	U,V,W coordinate (meters)		
- WEIGHT_SPECTRUM	weight per visibility		
- TIME	time in MJD (seconds)		
- ANTENNA1	first station of baseline (index in ANTENNA subtable)		
- ANTENNA2	second station of baseline		
- Several more columns			
- Subtables contain meta data			
- ANTENNA			
- NAME, POSITION,			
- FIELD			
- PHASE_DIR,			
— — —	- SPECTRAL_WINDOW		
- REF_FREQUENCY, - Several more subtables	CHAN_FREQ, CHAN_WIDTH,		
Show MS info, columns and subtables:			
showtable in=a	a.ms dm=f		
msoverview in-	=a.ms # msoverview -h for help		
See http://casacore.github	.io/casacore-notes/229.html for details		

showtable in=a.ms dm=f



Structure of table /Users/diepen/data/LOFAR L33277 SB010 uv.MS ----- Measurement Set 10353 rows, 23 columns (using 1 data managers) TIVW double shape=[3] unit=[m,m,m] measure=uvw,J2000 directly stored FLAG Bool ndim=2 ndim=3 FLAG CATEGORY Bool WEIGHT float ndim=1 float ndim=1 SIGMA ANTENNA1 Int scalar ANTENNA2 Int scalar Int scalar ARRAY ID DATA DESC ID Int scalar scalar unit=[s] EXPOSURE double FEED1 Int scalar FEED2 Int scalar FIELD ID Int scalar FLAG_ROW Bool scalar INTERVAL double scalar unit=[s] OBSERVATION ID Int scalar scalar PROCESSOR ID Int SCAN NUMBER Tnt scalar STATE ID scalar Int TIME double scalar unit=[s] measure=epoch,UTC TIME CENTROID double scalar unit=[s] measure=epoch,UTC Complex ndim=2 DATA WEIGHT SPECTRUM float ndim=2 SubTables: /Users/diepen/data/LOFAR L33277 SB010 uv.MS/ANTENNA /Users/diepen/data/LOFAR L33277 SB010 uv.MS/DATA DESCRIPTION /Users/diepen/data/LOFAR L33277 SB010 uv.MS/FEED /Users/diepen/data/LOFAR L33277 SB010 uv.MS/FLAG CMD /Users/diepen/data/LOFAR L33277 SB010 uv.MS/FIELD /Users/diepen/data/LOFAR L33277 SB010 uv.MS/HISTORY /Users/diepen/data/LOFAR L33277 SB010 uv.MS/OBSERVATION /Users/diepen/data/LOFAR L33277 SB010 uv.MS/POINTING /Users/diepen/data/LOFAR_L33277_SB010_uv.MS/POLARIZATION Dwingeloo, 25-Aug-2015 LOFAR synthesis/data/handling.77aQ10_uv.MS/PROCESSOR - 4 -/Users/diepen/data/LOFAR L33277 SB010 uv.MS/SPECTRAL WINDOW /Users/diepen/data/LOFAR L33277 SB010 uv.MS/STATE /HEORE/dionon/data/LOERE T22277 CON10 MC/LOERE CHAMTON

Some basic examples



```
taql 'select from a.ms where ANTENNA1!=ANTENNA2 giving cross.ms'
     tagl 'select from a.ms where ANTENNA1!=ANTENNA2 giving cross.ms as plain'
Select the cross-correlations and store the result in another table.
The first one as a RefTable (takes < 1 second); the second one makes a true copy (much slower).
    taql 'SELECT * FROM my.ms OFFSET 1e20'
Do an empty selection to show the column names only.
    taql 'select INTERVAL from my.ms limit 1'
Show the integration time (is constant in a LOFAR MeasurementSet)
    select result of 1 rows
1 selected columns: INTERVAL
4.00556
    taql 'select from my.ms orderby unique TIME'
    taql 'select from my.ms orderby unique ANTENNA1, ANTENNA2'
    tagl 'select gcount() from my.ms'
Show the number of time slots, baselines, and total number of rows.
    taql 'select ANTENNA1, ANTENNA2, UVW from my.ms
           orderby descending sumsqr(UVW[:2] limit 1'
Show the stations forming the longest baseline. Later we'll see how to get the names of the stations.
    taql 'insert into a.ms/STATE set SIG=True, REF=False, CAL=0, LOAD=0,
             SUB SCAN=0, OBS MODE="", FLAG ROW=False
Add a row to the STATE subtable and write the given column values in it.
```

TaQL commands



SQL-like data selection and manipulation can often replace a python script

- SELECT
 - select columns, select rows, sort rows
- UPDATE
 - update data in one or more rows and columns
- INSERT
 - add rows and fill them
- DELETE
 - delete rows
- CREATE TABLE
 - create a new table
- CALC
 - calculate an expression, possibly using table data
- COUNT
 - count number of rows for table subsets (e.g., per baseline)

SELECT



 Selects rows and /or columns and creates a new table result is normally a so-called RefTable (references the selection in the original table) Most parts are executed in the order given below, but SELECT after GROUPBY 				
SELECT columns				
FROM tables	columns or expressions to select (default all)			
	the input table(s) to use			
WHERE expression	which rows to select (default all); must result in bool scalar			
GROUPBY columns	scalar columns or expressions to group and aggregate on			
HAVING expression				
ORDERBY columns	groups to select			
LIMIT N	sort the result on scalar columns or expressions			
	maximum number of result rows (default all) (<0 is from end)			
OFFSET M	skip first M result rows (default 0); useful with ORDERBY (<0 is from end)			
GIVING table				
	persistent output table (default none)			

TaQL functionality



- Support of sets/arrays and many array functions
- Support of glob patterns and regex
- Support of units
- Support of date/time (UTC)
- Support of cone search
- Support of user defined functions (measure handling)
- Advanced interval support
- Support of nested queries
- Aggregation (GROUPBY, HAVING)
- Limited joins (only implicit equi-join on rownumber or id)
- Case-insensitive (except column names and string constants)

Looks a bit overwhelming, but simple selections can be expressed simply, especially in pyrap.tables May require careful thinking how to express a query Some SQL knowledge is advantageous See http://casacore.github.io/casacore-notes/199.html

Where can TaQL be used?



- in C++ casacore using function tableCommand
- in pyrap using function taql
 - indirectly in functions t.query, t.sort, t.select, and t.calc
- on command line using the program taql
 - use 'taql -h' to see how to use it
 - can also be used interactively (with command recall)
- Most important commands:
 - select
 - update
 - insert
 - calc

TaQL styles



TaQL indexing can have different styles. The standard way resembles Python

- array indices and row numbers start counting at 0
- end is exclusive
 - [0:5] is [0,1,2,3,4]
 - [1:3:0.5] is [1., 1.5, 2., 2.5]
- array axes order is C-style (row major)
 - first axis varies slowest
 - e.g., DATA column in an MS has axes [freq,pol]

Opposite is the old Glish style.

Simple queries



Simple queries can be expressed simply using some *pyrap.tables* functions. After

```
t = pt.table('my.ms')
          t1 = t.query ('ANTENNA1 = ANTENNA2')  # select auto-correlations
results in
          t1 = taql('select from my.ms where ANTENNA1 = ANTENNA2')
  in fact, in:
          t1 = tagl('select from $1 where ANTENNA1 = ANTENNA2', t)
                                                        # sort in time
          t2 = t1.sort ('TIME')
results in
          t2 = taql('select from $1 orderby TIME', t1)
          t3 = t2.select ('ANTENNA1, ANTENNA2')  # select a few columns
results in
          t3 = taql('select ANTENNA1, ANTENNA2 from $1', t2)
Combine as:
          t3 = t.query ('ANTENNA1=ANTENNA2', sortlist='TIME', columns='ANTENNA1, ANTENNA2')
results in
          t3 = taql('select ANTENNA1, ANTENNA2 from $1 where ANTENNA1=ANTENNA2 orderby TIME', t)
```

- 11 -

Data types



bool T, F, True, or False ٠ int64 double also sexagesimal format • 12h34m56.78 12d34m56.78 or 12.34.56.78 dcomplex 1 + 2i (or 2j) NB. normal addition • 3*2+2i == 6+2i not 6+6i in single and/or double quotes string 'a"b'"a'c" means a″ba′c 10-Nov-2010/12:34:23.98 datetime perl-like m/CS.*/ p/CS*/ regex ۲

Both scalars and arrays of these data types (NOT an array of regex)

A table column or keyword can have any table data type

Operators



On scalar and/or array; in order of precedence:

+ - & ^ == != ~= !~	+ - // % > >= < <= = ONE BETWEEN EXIST LIKE	power unary operators; ~ is bitwise complement // is integer division result; % is modulo; 1./2=0.5 1.//2=0 + is also string concatenation bitwise and bitwise and bitwise xor bitwise or normal comparison about equal (near function with 1e-5 tolerance) rs pattern matching logical and logical or
Operator n	ames are case-insensi	tive. For SQL compliancy some operators have a synonym.
=	=	
! =	<>	
& &	AND	
	OR	
!	NOT	
^	XOR	

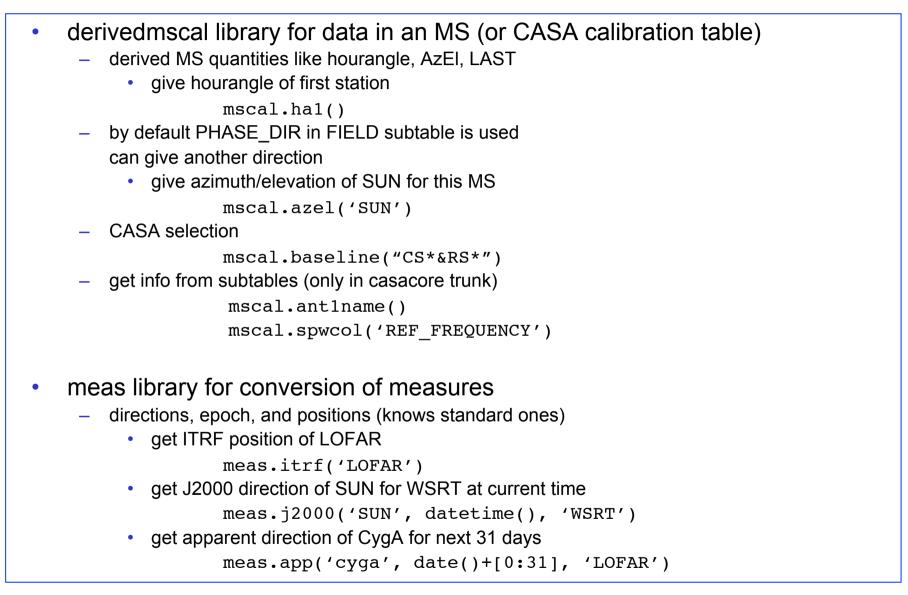
Functions



- > 160 standard functions
- Mathematical
 - pi, e, c, sqrt, sin, cos, asin, sinh, exp, pow, log, ...
- Comparison
 - near (operator ~=), nearabs, isnan, isinf, isfinite
- String, Regex
- Date/time
- Array reduction
 - mean, min, sumsqr, median, fractile, ...
 - plural form (mins, ...) for reduction of specific axes (e.g. per line)
 - sliding and boxed array functions
- Cone search
- Aggregation
 - gcount, gsum, gmean, gaggr, ...
- Miscellaneous
 - angdist, rand
 - iif (condition, val1, val2) (like ternary ?: in C)
 - Type conversion and string formatting
- User defined
 - xxx.func
 - are taken from shared library libcasa_xxx.so or libxxx.so

User defined functions





Units



Units are given implicitly or explicitly and converted if needed.

- A table column can have a unit
- Some functions result in a unit (e.g. *asin* results in unit rad)
- A sexagesimal constant has a unit (rad)
- A unit can be given after an expression. Conversion done if needed.
 Use quotes if a composite unit is used (e.g. 'km/s')

12 s	12 seconds	
12 s + 1 h	3612 seconds	
1 h + 12 s	1.00333 hour	
(174 lb)kg	78.9251 kg	(in case you use an American scale :-)
(1 'm/s')'km/h′	3.6 km/h	
12h34m56.78	3.29407 rad	
12 m < 1 km	True	
sin(45 deg)	0.707107	

These expressions can be given directly in taql program (assumes CALC if no command)

Regex



TaQL has rich support for pattern/regex matching (perl-like)

```
NAME ~ p/CS*/
                       match a glob pattern (as filenames in bash/csh)
NAME ~ p/CS*/i
                       same, but case-insensitive
NAME !~ p/CS*/
                       true if not matching the pattern
NAME ~ f/CS.*/
                       match an extended regular expression
NAME ~ m/CS/
                      true if part of NAME matches the regex (a la perl)
NAME ~ f/.*CS.*/
                       is the same
NAME like 'CS%'
                       SQL style pattern matching (also: not like)
NAME = PATTERN('CS*') glob pattern using a function (also !=)
NAME = REGEX('CS.*')
NAME = SQLPATTERN('CS\%')
Advanced
NAME ~ d/CS001/1 string distance (i.e., similarity)
```

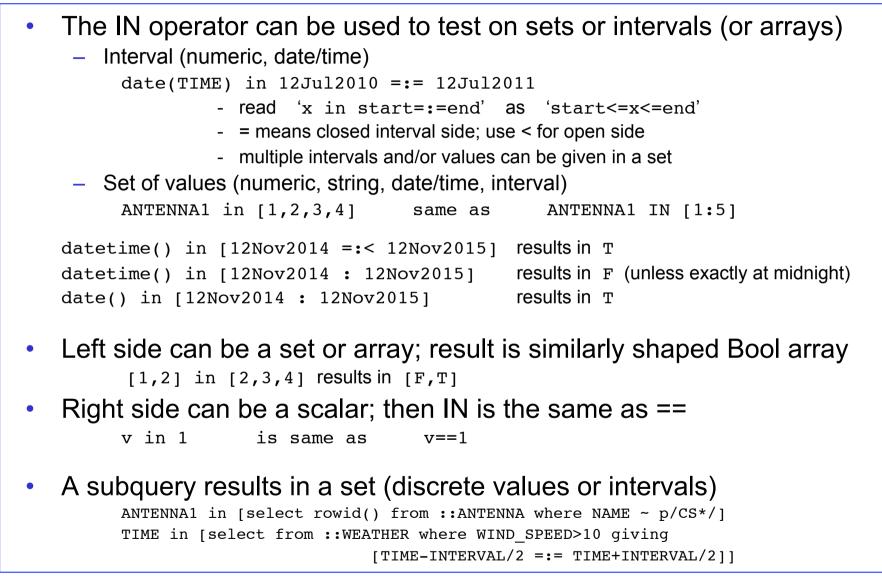
Arrays



•	Arrays can arise in several ways:		
	- a table column containing arrays		
	- a set results in a 1-dim array		
	[1:4] ['str1', 'str2']		
	- function array constructs an N-dim array		
	array([1:5], [3,4]) array with shape [3,4] filled with [1,2,3,4] in each row		
	array(DATA, product(shape(DATA)))) reshape to a vector		
•	Slicing is possible as in numpy (no negative values)		
	axes can be omitted (yields full axis)		
	DATA[,0] only take XX correlation from MS		
	DATA[::2,] take every other channel and all correlations		
•	Full array arithmetic (array-array, array-scalar, scalar-array)		
	 all operators, many functions (sin, etc.) available; they work element-wise 		
	 shapes have to match (no broadcasting like in numpy) 		
٠	Reduction functions (also partial for one or more axes)		
	 min, median, any, all, ntrue,, mins, medians, anys, alls, ntrues, 		
•	Sliding (running) functions		
	 runningmin, runningmedian, 		
•	Boxed functions		
	 boxedmin, boxedmedian, 		

Sets and intervals





Grouping and Aggregation



- Use to get aggregated information per group (e.g., baseline) groupby ANTENNA1, ANTENNA2
- Normally used with aggregation functions in SELECT

select gntrue(FLAG) from my.ms
select gntrue(FLAG) from my.ms groupby TIME

- Standard aggregation functions (all start with g)
 - gmin, gmax, gsum, gsumsqr, gmean, gstddev, gmedian, gfractile, gntrue, ...
- Special aggregation functions

gfirst, glastvalue in first/last row in groupgaggrconcatenate all rows of group into single arraymaxs(gaggr(abs(DATA)), 0,1)max amplitude per pol

• HAVING can be used to select groups

select gntrue(FLAG) as NT from my.ms groupby TIME having NT>0
Only selects groups where flags are set

UPDATE



Updates one or more columns in a table for each matching row **UPDATE** table The table to update SET column=expression, column=expression, ... The columns to update and their new values If column contains an array, a slice can be given A scalar can be assigned to an array (fills entire array) FROM table(s) Possible other input tables to use WHERE expression Which rows to update (default all); must result in bool scalar **ORDERBY** columns Sort scalar columns or expressions (default no sorting) LIMIT N Maximum number of matching rows to update (default all) OFFSET M Skip first M matching rows (default 0); useful with ORDERBY For example, set entire FLAG column to False UPDATE your.ms SET FLAG=False

CALC



Calculates an expression; if a table is given, it is calculated for each matching row

CALC expression				
the expression to calculate				
FROM tables the input table(s) to use (default none)				
WHERE expression				
which rows to use (default all); must result in bool scalar				
ORDERBY columns				
sort scalar columns or expressions (default no sorting) LIMIT N				
maximum number of matching rows to update (default all)				
OFFSET M				
skip first M matching rows (default 0); useful with ORDERBY				
For example:				
CALC ctod(TIME) from your.MS orderby unique TIME # format (unique) times				

Pretty printing using TaQL



```
By default program tagl pretty prints times
•
          taql 'select TIME from ~/GER1.MS orderby unique TIME limit 2'
              select result of 2 rows
          1 selected columns: TIME
          28-May-2001/02:26:55.000
          28-May-2001/02:27:05.000
   and positions
•
          taql 'select NAME, POSITION from ~/GER1.MS::ANTENNA'
              select result of 14 rows
          2 selected columns: NAME POSITION
          RT0
                  [3.82876e+06 m, 442449 m, 5.06492e+06 m]
          RT1
                  [3.82875e+06 m, 442592 m, 5.06492e+06 m]
   In python script use function ctod
•
          # Pretty print TIME like '2001/05/28/02:29:45.000'
          # Note the use of $t1 in the TaQL command;
          # The function tagl substitutes python variables given by $varname
          t1 = t.sort ('unique desc TIME', limit=18)
          pt.taql('calc ctod([select TIME from $t1])')
          # or by passing the times as a python variable; need to tell unit is s
          times = t1.getcol('TIME')
          pt.taql('calc ctod($times s)')
          # or the best way (cdatetime is a synonym for ctod)
          t1.calc ('cdatetime(TIME)')
```

Pretty printing using str(ing)



The function *str* can be used

- Converts values to strings
- Optional C-style format string or C++ width.prec
- Can also format date/time and angle using 'or-ed' format strings (as defined in class *MVTime*) and optional width

```
str(2rad, 'angle')  # +114.35.30
str(2rad, 'angle|10')  # +114.35.29.6125
str(4Mar1953, 'DMY|DAY|NO_TIME')  # Wed-04-Mar-1953
str(1+2i, '%f + %fj')  # 1.000000 + 2.000000j
str(123, 'value=%08d')  # value=00000123
str('abcdef', 4)  # abcd
```

Selection examples



```
Select all cross-correlations and save result (as a RefTable)
    select from your.ms where ANTENNA1 != ANTENNA2 giving your cross.ms
Select all cross-correlations and save result (as a PlainTable, thus deep copy); much slower
    select from your.ms where ANTENNA1 != ANTENNA2 giving your cross.ms as plain
Select all rows where ROW FLAG does not match FLAG
    select from your.ms where ROW FLAG != all(FLAG)
Select all rows where some, but not all correlations in a channel are flagged.
Note: ntrues determines #flags per channel; shape(FLAG) gives [nchan,ncorr]; true result if true for any channel
    select from your.ms where any(ntrues(FLAG,1) in [1:shape(FLAG)[1]])
Select some specific baselines (2-3, 2-4, 4-5, and 5-6)
Note: for a row containing e.g. baseline 2-5 you get [TTFF]&&[FFTF] → [FFFF]
    select from your.ms where any(ANTENNA1=[2,2,4,5] && ANTENNA2=[3,4,5,6])
Get the age (in days); could be used to test if an observation is old enough (note: date() is today)
Note use of :: in subtable name
    calc date() - 4Mar1953
    calc 20Jun2011 - date(TIME RANGE[0]) from your.ms::OBSERVATION
Get unique times
    select from my.ms orderby unique TIME
    select unique TIME from my.ms
```

More selection examples



```
Select baselines (auto and cross) between LOFAR core and remote stations only
Note: this can also be achieved using the program msselect !!
    select from my.ms where mscal.baseline('[CR]S*&&') # CASA selection syntax
    select from my.ms where all([ANTENNA1,ANTENNA2] in
                                    [select rowid() from ::ANTENNA where NAME ~ p/[CR]S*/])
Select baselines containing international stations
    select from my.ms where mscal.baseline('^[CR]S*&&*')
Regression test; check if DATA column is as expected (NaNs and rounding errors are possible)
Note: t1.DATA ~= t2.DATA
                               is the same as
                                               near(t1.DATA, t2.DATA, 1e-5)
    select from test.ms t1, ref.ms t2 where
         not all((isnan(t1.DATA) and isnan(t2.DATA)) or t1.DATA ~= t2.DATA)
```

More selection examples



Show short baselines (< 200 m) with the antenna names (with some kind of join using indexing)

```
tagl 'select ANTENNA1, ANTENNA2, sqrt(sumsqr(UVW[:2])),
           [select NAME from ::ANTENNA][ANTENNA1] as ANTNAME1,
           [select NAME from ::ANTENNA][ANTENNA2] as ANTNAME2
    from ~/DATA/GER.MS where sumsqr(UVW[:2]) < 200*200 orderby unique ANTENNA1, ANTENNA2'
5 selected columns: ANTENNA1 ANTENNA2 Col 3 ANTNAME1 ANTNAME2
    1
0
          73.9286
                    RT0
                              RT1
0
   2
         147.859
                    RT0
                              RT2
      73.8642
1
  2
                    RT1
                              RT2
      147.723
1
    3
                              RT3
                    RT1
With newest version of casacore (use Loflm)
  taql 'select ANTENNA1, ANTENNA2, sqrt(sumsqr(UVW[:2])),
           mscal.antlname() as ANTNAME1,
           mscal.ant2name() as ANTNAME2
    from ~/DATA/GER.MS where sumsqr(UVW[:2]) < 200*200 orderby unique ANTENNA1, ANTENNA2'
```

Update examples



Clear all flags in a MeasurementSet or set the flag if corresponding DATA is invalid update your.ms set FLAG=False, ROW FLAG=False update your.ms set FLAG=!isfinite(DATA), ROW FLAG=all(FLAG) # uses new FLAG values!! Update the MOUNT in the ANTENNA table update your.ms/ANTENNA set MOUNT='X-Y' Put CORRECTED DATA into the DATA column update my.ms set DATA = CORRECTED DATA Put CORRECTED DATA from that MS into DATA column of this MS It requires that both tables have the same number of rows update this.ms, that.ms t2 set DATA = t2.CORRECTED DATA Subtract background noise from an image using the median in a 51x51 box around each pixel (updates the image, so one should have made a copy of my.img) update my.img set map = map - runningmedian(map, 25, 25) Flag XX data based on a simple median filter (per row); if set, current flag is kept update my.ms set FLAG[,0]=FLAG[,0] || amplitude(DATA[,0]) > 3*median(amplitude(DATA[,0])) where isdefined(DATA) Add a line to the HISTORY table of a MeasurementSet (converts the MJD time automatically to sec) insert into my.ms/HISTORY (TIME, MESSAGE) values (mjd(), "historystring")

Calculation examples



```
Count all flags in an MS (uses nested query)
    CALC sum([select ntrue(FLAG) from my.ms])
Get percentage of unflagged data in an MS
    CALC sum([select nfalse(FLAG) from my.ms]) * 100. /
          sum([select nelements(FLAG) from my.ms])
Get the hourangle of the first station (creates a new PlainTable, not RefTable because of expression in columns)
    SELECT TIME, ANTENNA1, ANTENNA2, mscal.hal() as HA1 from my.ms giving newtable
The same, but return it a as an array
    CALC mscal.hal() from my.ms orderby unique TIME, ANTENNA1
Angular distance between observation's phase reference direction(s) and a given direction
Do the same for CyqA
    CALC angdist([-3h45m34.95, 10d12m42.5], PHASE DIR[0,]) FROM your.ms/FIELD
    CALC angdist(meas.j2000('CygA'), PHASE DIR[0,]) FROM your.ms/FIELD
Angular distance (in radians) between apparent positions of sun and moon at LOFAR core in coming month
(sun and moon are close on 5-Dec, so it'll be dark for poor Sinterklaas and Zwarte Piet)
     CALC angdist(meas.app('SUN', 26Nov2013+[0:31], 'LOFAR'),
                    meas.app('MOON',26Nov2013+[0:31], 'LOFAR'))
    [1.53393, 1.33603, 1.13221, 0.921343, 0.702766, 0.476634, 0.244891, 0.0500374, 0.24891, 0.490751,
    0.732487, 0.970688, 1.20361, 1.43043, 1.65101, 1.86568, 2.07504, 2.27972, 2.48028, 2.67696, 2.86884,
    3.04371, 3.00427, 2.82496, 2.63905, 2.45208, 2.26415, 2.07456, 1.88225, 1.68595, 1.48433]
```

uvflux example



The Miriad program uvflux estimates the source I flux density and its standard deviation at the phase center without having to make an image.

A single, not too complicated TaQL command (courtesy Dijkema, Heald) provides the same functionality.

```
select gstddev(SUMMED) as STDVALS,
4.
             gmean(SUMMED) as MEANVALS,
4.
             gcount(SUMMED) as NVALS
4.
      from (select gmean(
3.
1.
                    sum(iif(FLAG[,0:4:3], 0, abs(DATA[,0:4:3]))) / nfalse(FLAG[,::3])
3.
                    ) as SUMMED
            from ~/data/GER.MS
            where mscal.baseline('5km~10km') && !all(FLAG)
2.
3.
            groupby TIME)
```

A subquery is made to get the mean I flux (= $0.5^{(XX+YY)}$) per time slot in the following way.

1. It first gets the mean of all channels for each baseline. Note that it uses sum/n to ignore flagged visibilities. The iif function tells to use 0 for them. Also note that XX is the 1st, YY the 4th polarisation, hence [0:4:3] (or [::3]) indexes these polarisations.

Once masked arrays are supported by TaQL, it could be written as: mean(DATA[,::3][FLAG[,::3]])

- 2. It only uses the baselines with lengths between 5 and 10 km where not all visibilities are flagged.
- 3. Thereafter the average flux per time slot is determined in the subquery using the gmean aggregation and GROUPBY functionality. The result is called SUMMED.
- 4. Finally, the outer query uses aggregate functions to calculate the overall mean, standard deviation, and number of time slots. The final result is a table with 1 row and 3 columns.

Check StationAdder result



The StationAdder step in DPPP forms a new station ST001. It forms new baselines ST001 with all non-core stations by adding data of the baselines of non-core stations with core stations. Per timeslot we want to check for each new baseline if the resulting DATA and UVW are as expected.

NewData = sum(OldData * Weight) / sum(Weight)# only use unflagged data pointsNewUVW = sum(UVW * SumWeight)) / sum(SumWeight))# SumWeight is sum of weights of all channels

1. In a LOFAR MS non-core/core baselines have ANTENNA1 as core and ANTENNA2 as non-core!!

2. Per time slot and non-core station we have to combine (group) the data -> use GROUPBY

select from a.ms groupby TIME,ANTENNA2

3. Only use unflagged data -> use 0 for flagged data points

iif(FLAG, 0, DATA)

```
upcoming version of TaQL will support masked arrays -> DATA[FLAG]
```

4. Use the non-core/core baselines -> selection of such baselines most easily done with CASA syntax

```
where mscal.baseline("^CS*&CS*")
```

5. Calculate new data by aggregating the data per group and summing over the first axis (i.e., baselines)

```
sums(gaggr(DATA*WEIGHT_SPECTRUM), 0) / sums(gaggr(WEIGHT_SPECTRUM), 0)
```

6. Combine it all, write the result in an output table, and call the output column NDATA

select from avg.data t1, [select from a.ms where mscal.baseline("ST001&*"] t2

where not all (t1.NDATA ~= t2.DATA)

Check StationAdder's UVW



The following tagl command calculates the average UVW per new baseline taking the data flags into account. It was used by Leah Morabito to visually check DPPP's StationAdder results. If the results are stored in a table, another TaQL command could be used to compare with the StationAdder output. It resembles the previous uvflux example. It uses a subquery to calculate intermediate results. It also resembles numpy with its whole and partial array operations. This guery takes advantage of the knowledge that in a LOFAR MS the core stations are in ANTENNA1 for the baselines between core and non-core baselines, thus it can sum over ANTENNA2. select ctod(TIME), ANTENNA2, 4. sums(gaggr(UVW*SUMW),0) / gsum(SUMW) as UVW 1. from [select TIME,UVW,ANTENNA2, sum(iif(FLAG,0,WEIGHT SPECTRUM)) as SUMW 3. 1. from a.ms 2. where mscal.baseline("^CS*&CS*")] 4. groupby TIME, ANTENNA2 1. The subquery selects the required columns. 2. It only uses the rows with cross-correlation baselines between non-core and core (CASA syntax). 3. As part of the subquery output it calculates SUMW, the sum of the weights where the data are not flagged. Note: in the next release of casacore one can probably use masked arrays like: sum(WEIGHT SPECTRUM[FLAG]) 4. The main guery aggregates the subguery output to a single array per TIME and ANTENNA2. It sums each array over the first axis, resulting in the weighted sum for U, V, and W. Finally it is divided by the summed SUMW to get the weighted averaged U, V, and W per TIME and ANTENNA2. Dwingeloo, 25-Aug-2015 LOFAR synthesis data handling: TaQL

Advanced examples 1



Check if an MS is in non-descending time order and check for missing time slots. It is done by comparing the table subsets of row 0..n-1 and row 1..n (created by a nested query in the FROM). Note the use of LIMIT -1 to denote the one but last row (a la python indexing).

```
select t1.TIME, t1.TIME-t2.TIME as STEP from
  [select from ~/L33277_SAP000_SB000_uv.MS limit -1] t1,
   [select from ~/L33277_SAP000_SB000_uv.MS offset 1] t2
  where int((t2.TIME-t1.TIME)/t1.EXPOSURE+0.0001) not in [0,1]
```

Check if the QUALITY subtables generated by NDPPP and rficonsole are the same, except for the baselines that were flagged already (because rficonsole counts them, but NDPPP does not). Note that t3 contains the baselines with some flagged data in the original MS.

```
select from ndppp.ms/QUALITY_BASELINE_STATISTIC t1,
    rficonsole.ms/QUALITY_BASELINE_STATISTIC t2,
    [select unique ANTENNA1,ANTENNA2 from original.ms where any(FLAG)] t3
where any(t1.VALUE != t2.VALUE) and
!any(ANTENNA1 = [select ANTENNA1 from t3] && ANTENNA2 = [select ANTENNA2 from t3])
```

Advanced examples 2



Check if the demixing solutions in the old and new way are the same for CygA. The old one only contains CygA, the new one contains more, so a selection is needed.

```
taql 'select from instrumentold t1,
        [select from instrument3 where NAMEID in
      [select rowid() from ::NAMES where NAME~m/CygA/]] t2
   where !all(t1.VALUES ~= t2.VALUES)'
```

```
select result of 0 rows
```

Note that above assumes that the table orders are the same. If not,

orderby STARTX, STARTY should be used for both tables to make them the same.

```
taql 'select from [select from instrumentold orderby STARTX,STARTY] t1,
        [select from instrument3 where NAMEID in
        [select rowid() from ::NAMES where NAME~m/CygA/] orderby STARTX,STARTY] t2
    where !all(t1.VALUES ~= t2.VALUES)'
```

Advanced examples 3



```
Swap columns ANTENNA1 and ANTENNA2 in a MeasurementSet.

The problem is that TaQL updates in place, thus as soon as one ANTENNA column is set, its original values are lost.

The first solution works fine (think about it), but is some kind of a hack.

update my.ms set ANTENNA1 = ANTENNA1+ANTENNA2,

ANTENNA2 = ANTENNA1-ANTENNA2, ANTENNA1 = ANTENNA1-ANTENNA2

The following solution is neater.

It holds the original values of ANTENNA1 in a temporary table in memory.

update my.ms, [select ANTENNA1 from my.ms giving as memory] as orig
```

set ANTENNA1 = ANTENNA2, ANTENNA2 = orig.ANTENNA1