H2020-Astronomy ESFRI and Research Infrastructure Cluster (Grant Agreement number: 653477).



## Task-based distributed processing for radio-interferometric imaging with CASA

BOJAN NIKOLIC

### 2<sup>nd</sup> ASTERICS-OBELICS Workshop

16-19 October 2017, Barcelona, Spain.



## Motivation

# The "embarrassing" parallelism in Radio Astronomy



Each observation is typically an hour of time

8700 hours in a year => almost 10000 way parallelism if we are patient

## That is enough to scale the SKA data reduction to a single node! But:

UNIVERSITY OF

CAMBRIDGE

- Sometimes can't wait! Calibration, fast transients
- Very large storage requirements about 1MW of spinning drives



## Traditional approaches

Use a single, high specification, computer and be patient Ideally with an automated pipeline script,

But sometimes :

#### The human pipeline: distributed data reduction for ALMA

Scott L. Schnee<sup>\*a</sup>, Crystal Brogan<sup>a</sup>, Daniel Espada<sup>b,c,d</sup>, Elizabeth Humphries<sup>e</sup>, Shinya Komugi<sup>c</sup>, Dirk Petry<sup>e</sup>, Baltasar Vila-Vilaro<sup>b</sup>, and Eric Villard<sup>b</sup>

 <sup>a</sup>National Radio Astronomy Observatory, 520 Edgemont Rd, Charlottesville, VA, U.S.A.;
<sup>b</sup>Joint ALMA Observatory, Alonso de Cordova 3107, Vitacura, 763-0355, Santiago, Chile;
<sup>c</sup>National Astronomical Observatory of Japan (NAOJ), 2-21-1 Osawa, Mitaka, 181-8588, Tokyo, Japan; <sup>d</sup>The Graduate University for Advanced Studies (SOKENDAI), 2-21-1

# Objectives for task based parallelisation

### Time-to-solution

- Non observation-parallel use casses
- Free-up expensive fast storage space as quickly as possible
- Feedback on quality of observations quicker
- Scientific results sooner
- Need fewer People

#### Distributed memory

- Avoid the expense and scaling limitation of large shared memory cachecoherent machine
- Typical HPC nodes have 64/128 GB memory while increasing trend to 1TB+ RAM machines

## Efficient use of I/O bandwidth

- Overlap of I/O and processing inside an application difficult to arrange
- Many simultaneous tasks allow a queue of I/O requests to built up





## Tasks



### Tasks

- Take arguments and produce a result
- 2. Can decompose into further tasks

What is the minimum constraint on ordering on tasks? How do we avoid over-specifying this ?



M. J. Flynn, "Some computer organizations and their effectiveness," IEEE Trans. Computers, vol. 21, no. 9, pp. 948–960, Sep. 1972.



### Tasks

- Take arguments and produce a result
- 2. Can decompose into further tasks

Does (lexical & runtime) program structure correspond to tasks?



M. J. Flynn, "Some computer organizations and their effectiveness," IEEE Trans. Computers, vol. 21, no. 9, pp. 948–960, Sep. 1972.

### Von Neumann Model

Instruction == Task

Ordering == Program Counter



### UNIVERSITY OF CAMBRIDGE

### Tasks

Precedence definition

### $f_i^l(x,\tau) \to (y,\tau^*)$

x : Input data

 $\tau$  : Input control variable, i.e., x valid for use in computation? y : Task result

 $\tau^*$ : Output control variable, i.e., is y valid?



# QCD: domain decomposition & fine grain global iteration



From Ukawa 2014:

https://indico.cern.ch/event/284433/contributions/1634880/attachments/525606/724848/ukawa.pdf



# QCD: domain decomposition & fine grain global iteration



#### From Ukawa 2014:

https://indico.cern.ch/event/284433/contributions/1634880/attachments/525606/724848/ukawa.pdf



## AIPS++ vs CASA

#### "TOOLS": AN OBJECT ORIENTED USER INTERFACE



#### "TASKS": A TASK-BASED INTERFACE





# CASA Script and equivalent task graph





## Architecture



## Drivers

Modular, readable, source code structure. Small increment from standard CASA usage

Short development time, low maintenance

Easy to use on standard HPC clusters and single nodes

Reasonable performance and scalability



## SWIFT/T

LANGUAGE

## EXECUTION ENGINE (TURBINE)

## Functional language with only I-Structure variables

- Every (lexical) variable can only be assigned to at most once
- This enables a compile time and run-time exact dataflow inference

### First class support for files and arrays

Simple, familiar looking, syntax (not Haskell!)

#### Simple but adequate type system

## Designed for HPC clusters with a shared filesystem (Lustre/GPFS) and MPI-3 libraries

High task issue rate

Fully dynamical load-balancing

## Capability for in-process tasks (including CASA!)

• No overhead of starting a process for each task

Integrates with cluster schedulers (SLURM, etc). Works out of box on typical cluster.

See Wilde (2011) for the SWIFT language; See Wozniak (2013) for description of SWIFT/T



## Module Architecture (simplified)





## RUN TIME ARCHITECTURE (COMPONENTS)



Drivers: Scalability, Load balancing, high task issue rate, low task issue overhead



## CASA to SWIFT/T





# SWIFT/T program and equivalent graph





## RESULTS



## MAIN RESULTS

Short development cycle

- Two weeks to full implementation (but I was already very familiar with CASA & SWIFT/T)
- Minimal need for interaction with upstream developers

#### Understandable program structure

- SWIFT if anything more understandable from scientific viewpoint than the underlying Python
- Accurate, reliable, modularisations



## Scalability





## More information

Memo with the details full source code



## Acknowledgement

H2020-Astronomy ESFRI and Research Infrastructure Cluster (Grant Agreement number: 653477).