# HPC-UGent pilot kickoff meeting

# `doduo` Tier-2 cluster

Oct 28th 2020

https://www.ugent.be/hpc/en/support/pilot/doduo

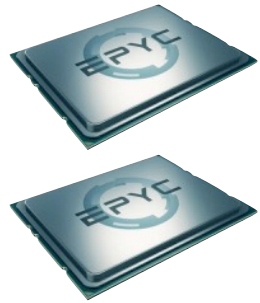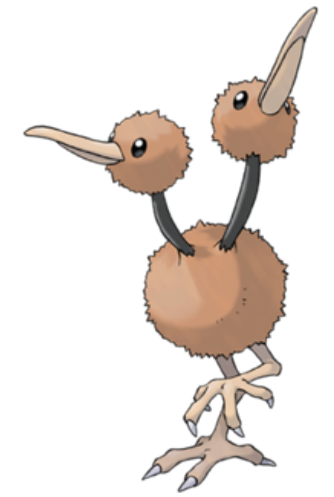hpc@ugent.be                                        http://ugent.be/hpc

# Pilot users for doduo

- members of `gpilot` user group (invitation only)

- experienced users of existing HPC-UGent Tier-2 infrastructure

- different research domains & applications

- mailing list: hpc-pilot-users@lists.ugent.be

  - used by HPC-UGent team to contact pilot users (status updates, etc.)

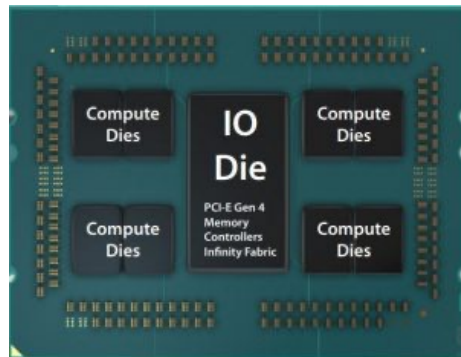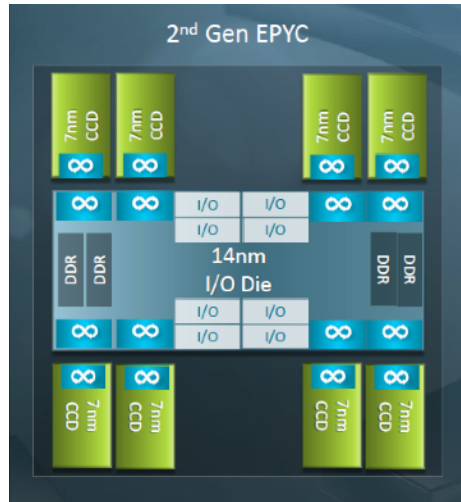  - can be used by pilot users to get in touch with each other

# Technical details for doduo

- **128 workernodes**, each with:

  - two 48-core AMD EPYC 7552 CPUs (AMD Rome)
    => **96 cores per node**

  - ~250GB usable RAM memory => ~2.5GB/core

  - ~180GB of local disk (SSD)

- **12,288 cores** in total

- **HDR-100 Infiniband** interconnect

- fast access to shared filesystems (GPFS)

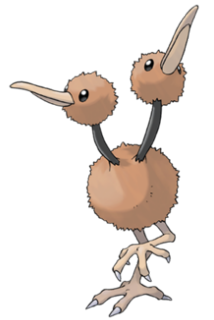- OS: **Red Hat Enterprise Linux 8.2 (RHEL8.2)**
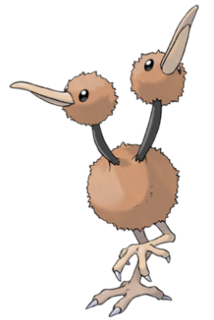
# Technical details: AMD Rome

- naming mess:
  - AMD EPYC (line of AMD processors)
  - AMD Rome (2nd generation of AMD EPYC)
  - AMD Zen2 (CPU microarchitecture for AMD Rome processors)

- 48-core AMD EPYC 7552 processor (2 per node in doduo)
  - 2.2GHz base clock, boost up to 3.3GHz
  - 512KB L2 cache per core
  - 192MB L3 cache (shared)
  - supports SSE4.2, AVX, AVX2, FMA
  - **does *not* support AVX512**

https://www.nextplatform.com/2019/08/15/a-deep-dive-into-amds-rome-epyc-architecture
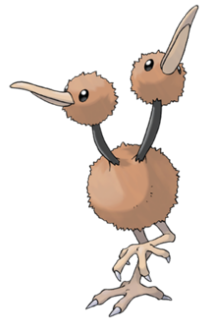
# Differences with existing Tier-2 clusters

- **AMD x86_64 processors** (vs Intel x86_64 processors)
  - complicates software installations (w.r.t. compilers/libraries to use)
  - may affect software that is sensitive to floating-point accuracy (VASP, CP2K, ...)

- **RHEL8.2** (vs CentOS Linux 7.8)
  - software built for doduo won't run on login nodes or other Tier-2 clusters! (GLIBC errors)

- **different wrappers for qsub, qstat, ... commands**
  - provided via new `jobcli` project, with Torque frontend + Slurm backend

- **96 cores** per node (vs 36 max.)
  - be careful with requesting full nodes, check scaling across cores first!

- **only recent toolchains**
  - foss: 2019b, 2020a
  - intel: 2019b (with newer impi version), 2020a (with *older* imkl version!)
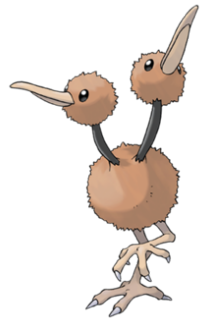
# Getting access to `doduo`

- submit jobs from HPC-UGent Tier-2 login nodes

  - only for members of `gpilot` user group!

- `module swap cluster/.doduo`

  hidden cluster module!

- if you compile any software yourself for `doduo`,
  **make sure you do that *on the workernodes* !**

  - login nodes: Intel Skylake (AVX2, AVX512) + CentOS 7

  - doduo: AMD Zen2 (AVX2 only) + RHEL8

GHENT
UNIVERSITY

# Scientific software on doduo

- request missing software installations via new form:

  https://www.ugent.be/hpc/en/support/software-installation-request

- **only with `2019b` or `2020a` toolchains** (or more recent)

  - recent compilers/libraries are required for RHEL8 / AMD Rome

  - strong preference for *latest* software versions

# Scientific software on `doduo`
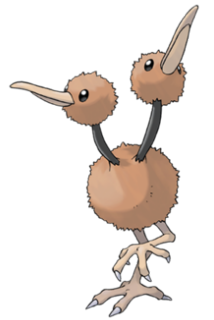
Currently available (see `module avail`):

CP2K (v7.1)          DIRAC          DL_POLY_Classic          Gaussian (g16_C.01)

GROMACS          LAMMPS          OpenFOAM          OpenMM

Python (+ SciPy-bundle)          SCons          VASP (v5.4.4)          VTK          yaff

Work in progress:

- Crystal17
- hanythingondemand (HoD)
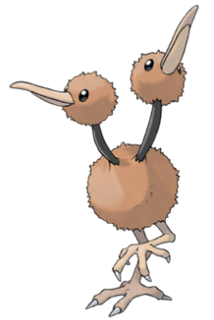- iPI
- RASPA (?)
- VASP 6

# mympirun

- new `mympirun` version (5.2.2)

- should work just like previous versions (with less warnings)

- make sure to always load latest version (don't specify a version)
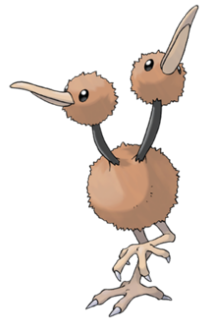
```
module load vsc-mympirun
```

- please report any problems ASAP via `hpc@ugent.be`

- future update will switch to different MPI startup mechanism
  (more on that later via hpc-pilot-users mailing list)
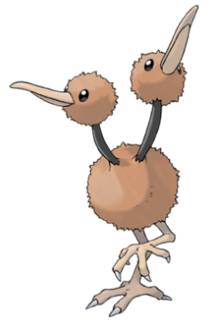
# Expectations from pilot users

- **testing** usage of new cluster & provided software

- **comparing** with existing Tier-2 (swalot, skitty) & Tier-1 clusters
    - re-run jobs => **validate results** + compare performance
    - try to make comparisons 'honest' (same or similar node/core count)
- **scaling tests** for parallel software (both intra-node and inter-node)
    - don't hold back, try large runs!

- **offload work from current Tier-2 clusters** (especially when golett is gone)

- **report back** findings (& problems) to `hpc@ugent.be`
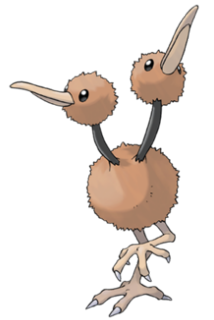
# Attention points (1/2)

- **be very critical** w.r.t. (scientific) results obtained during pilot phase

  - AMD Rome & Intel compilers/libraries may cause inaccuracies...

  - ~~double~~ triple-check your results!

  - report problems to `hpc@ugent.be` so we can mitigate if needed

- start with small experiments & re-running stuff you've run before

- gradually scale up, run new things when you're more confident

- **don't blindly use full nodes (96 cores each), check scaling first!**
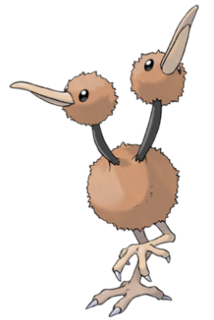
# Attention points (2/2)

- **pilot cluster should work like existing Tier-2 clusters**

  - there should be no need to change workflow/job scripts (other than `module load` statements)

  - *if you need to change something to get it to work (well), let us know!*

- **pilot clusters may become unavailable on (very) short notice**

  - down for maintenance to resolve problems or install updates

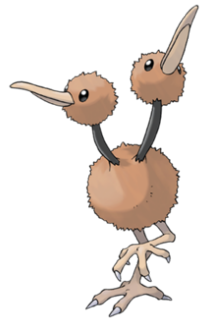  - unexpected downtime due to software/hardware/DC problems

# Known issues

- relatively low memory bandwidth per core

- aggressive power saving, may impact IB network performance

- minor issues with "job commands" (qsub, qstat & co)

  - doesn't work yet: `qdel all`, `qstat -t`

  - `qstat` is slow when there are lots of jobs

  - interactive jobs (`qsub -I`) start one shell per requested core

- software-specific issues

  - more failing tests in CP2K regression test (see installation log)

# Timeline (preliminary)

- *Oct 28th 2020 (today):*

  - access for pilot users (only `gpilot` members)

  - first 64 nodes available (part 1)

- all requested software available will be installed ASAP

- mid Nov'20: add part 2 (64 more nodes) + downtime part 1

- move from pilot to production: not before Feb'21

# Problems or questions?

- contact `hpc@ugent.be`

- make it clear in e-mail subject that it's related to pilot clusters

- provide clear problem description

  - what did you expect to work, what went wrong

  - mention relevant error messages, job IDs, etc.

  - mention location of output files in your account (please don't send them in attachment)

  - exact steps to reproduce the problem