



Introduction to research in Computer Science

François Trahay



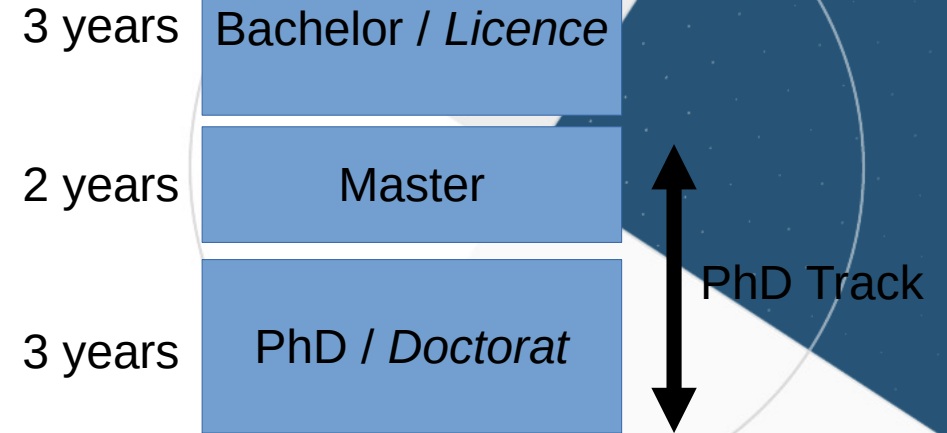
- What is a PhD ?
 - Grants
 - Finding a grant
- Publications
- Coping with the PhD
- Careers

What is a PhD ?

- 3 years work on a research topic
 - In a research lab / in a company
 - Work as an employee (salary : 1700 - 2200 € per month)

- Output of the PhD work
 - Research results (publications, software, patents, ...)
 - A PhD thesis manuscript
 - In CS : approx 100 pages manuscript
 - A PhD defense
 - 45 minutes presentation of the research work + many questions

-> A PhD degree

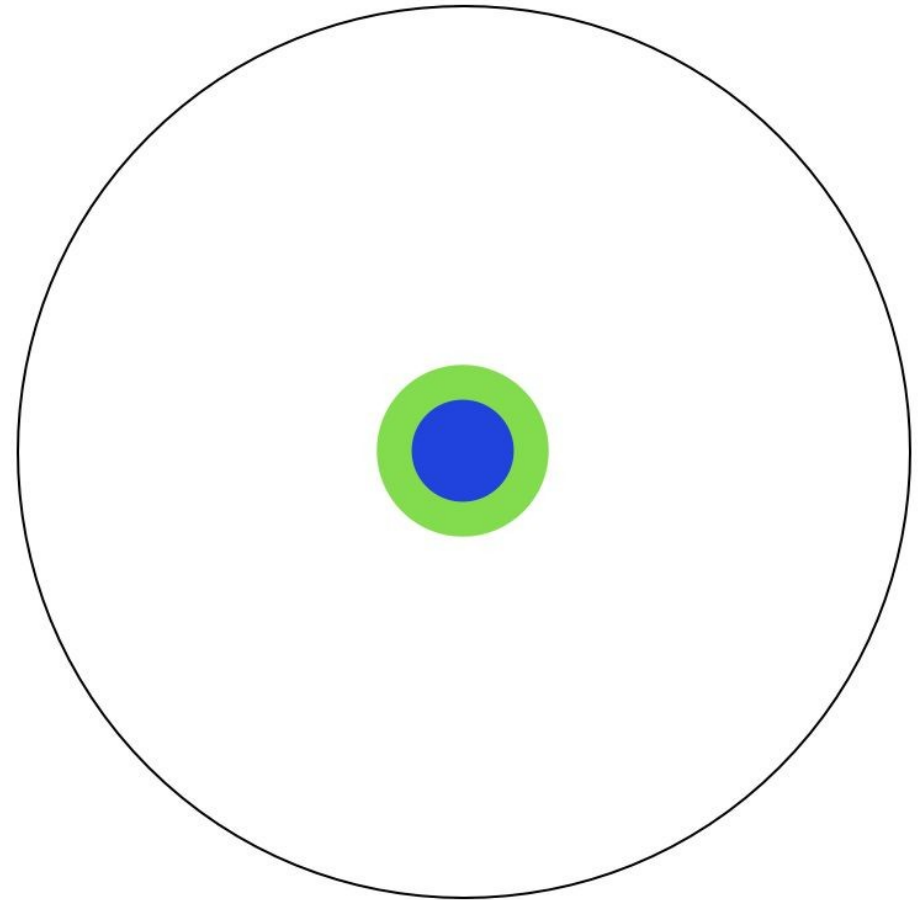


Why would you want to do that ?

- Doing research is fun
- You get the opportunity to work on what you want
- A thesis can open job opportunities
 - In the industry
 - In academia
- You help to advance science

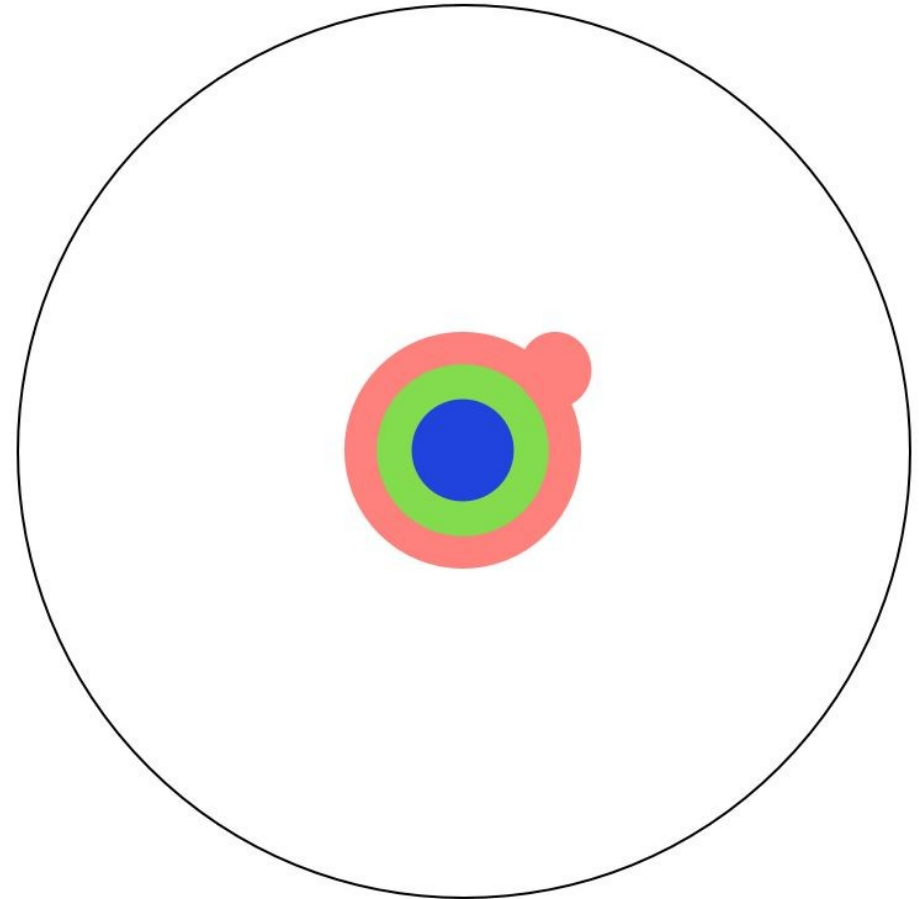
Human knowledge

Imagine a circle that contains all of human knowledge.
By the time you finish high school, you know a bit of ma



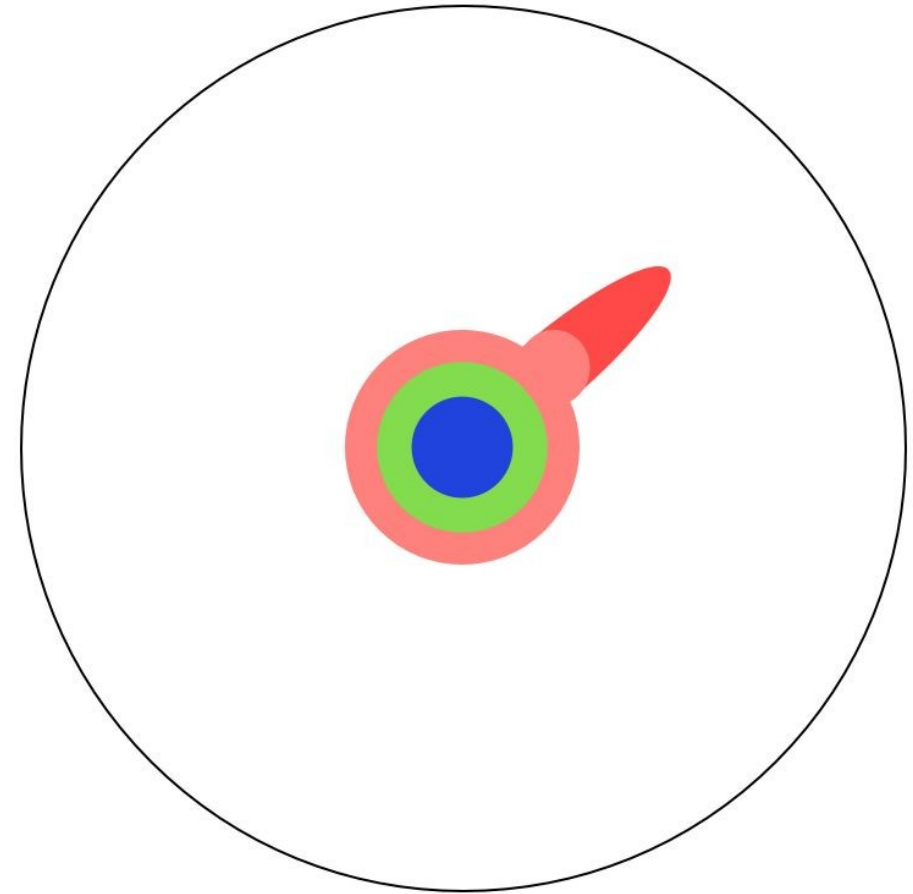
Human knowledge

With a bachelor's degree, you gain a specialty.



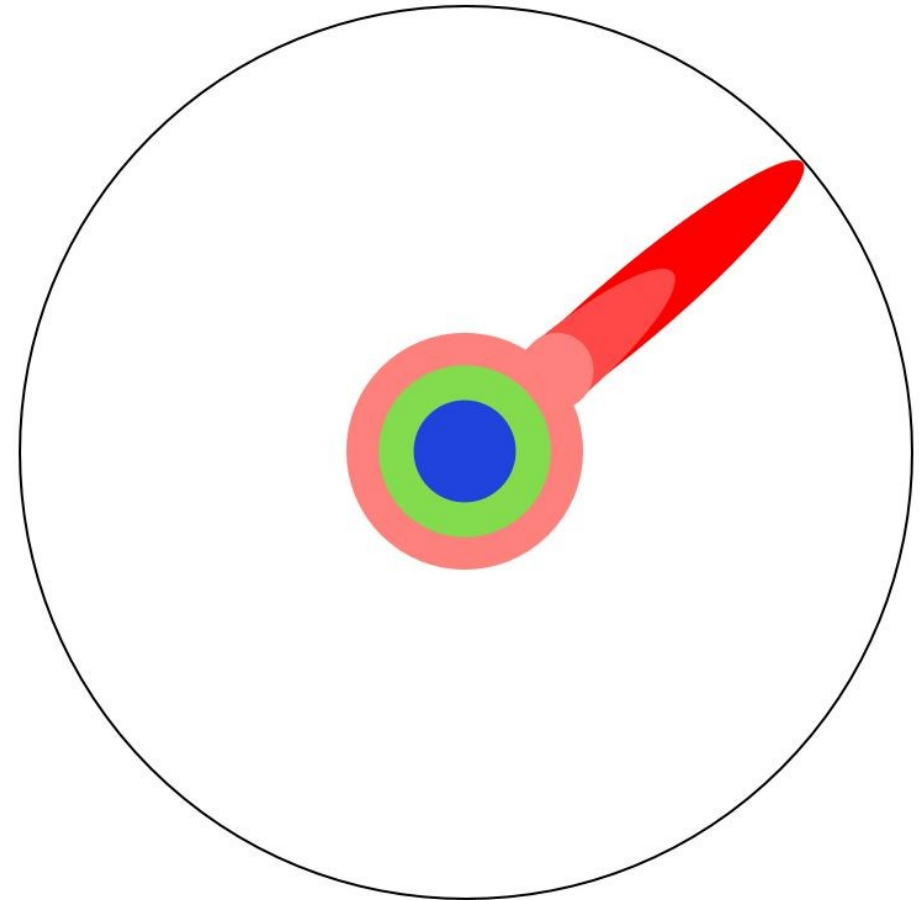
Human knowledge

A master's degree deepens that specialty.



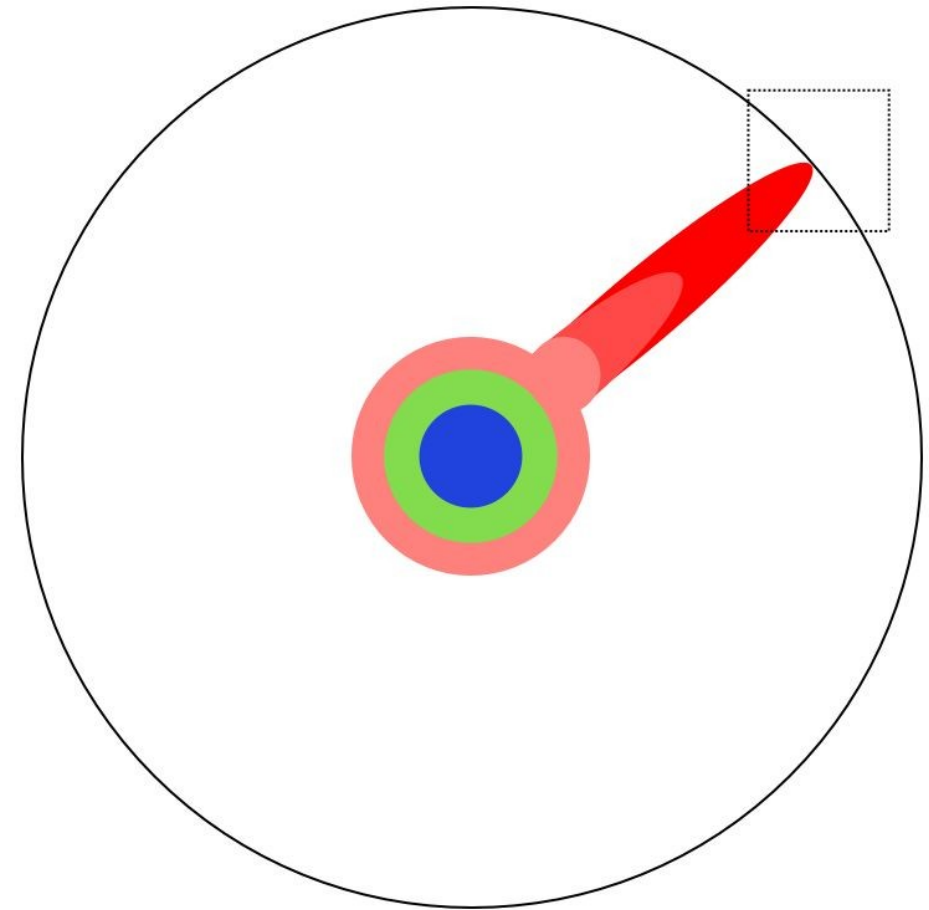
Human knowledge

Reading research papers takes you to the edge of human knowledge



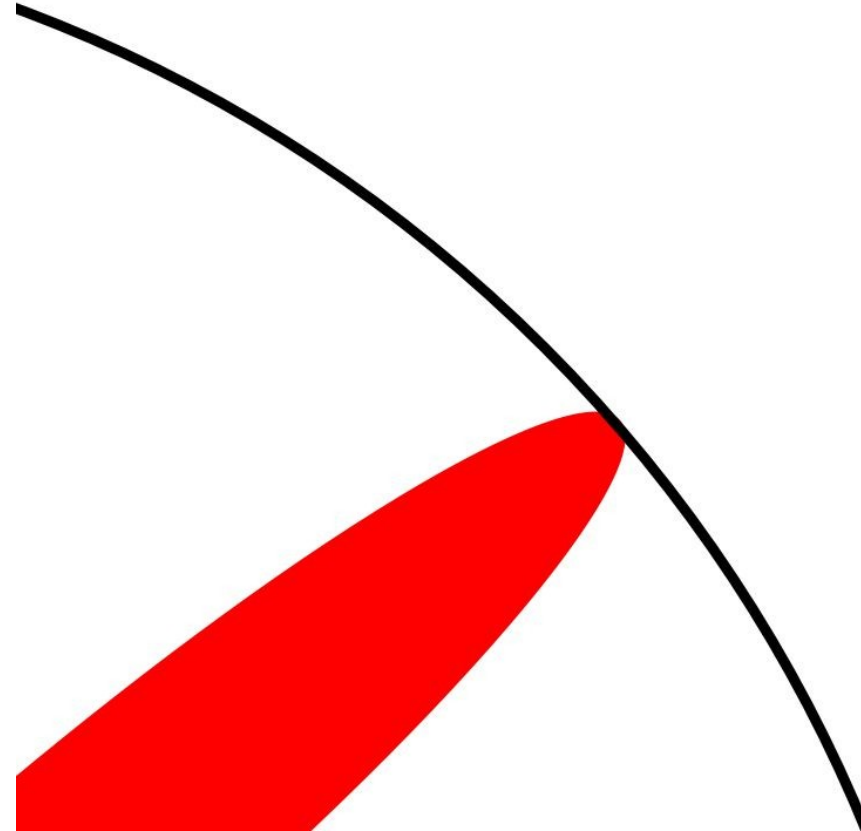
Human knowledge

Once you're at the boundary, you focus.



Human knowledge

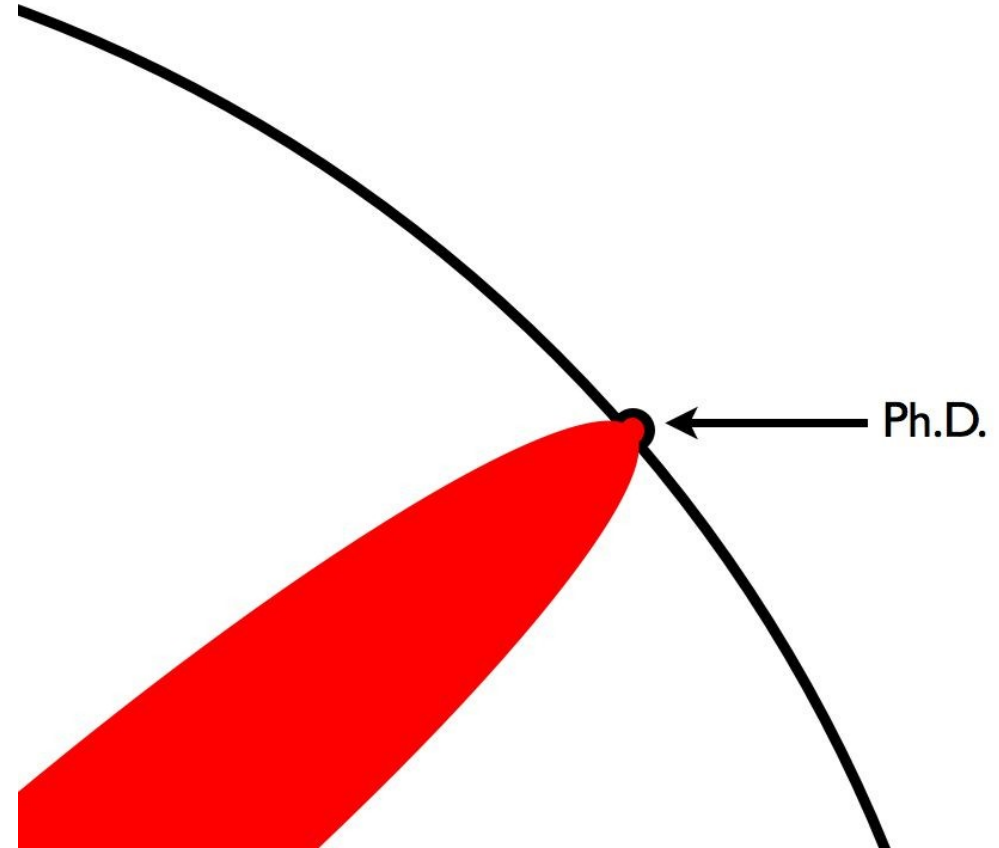
You push at the boundary for a few years.



Human knowledge

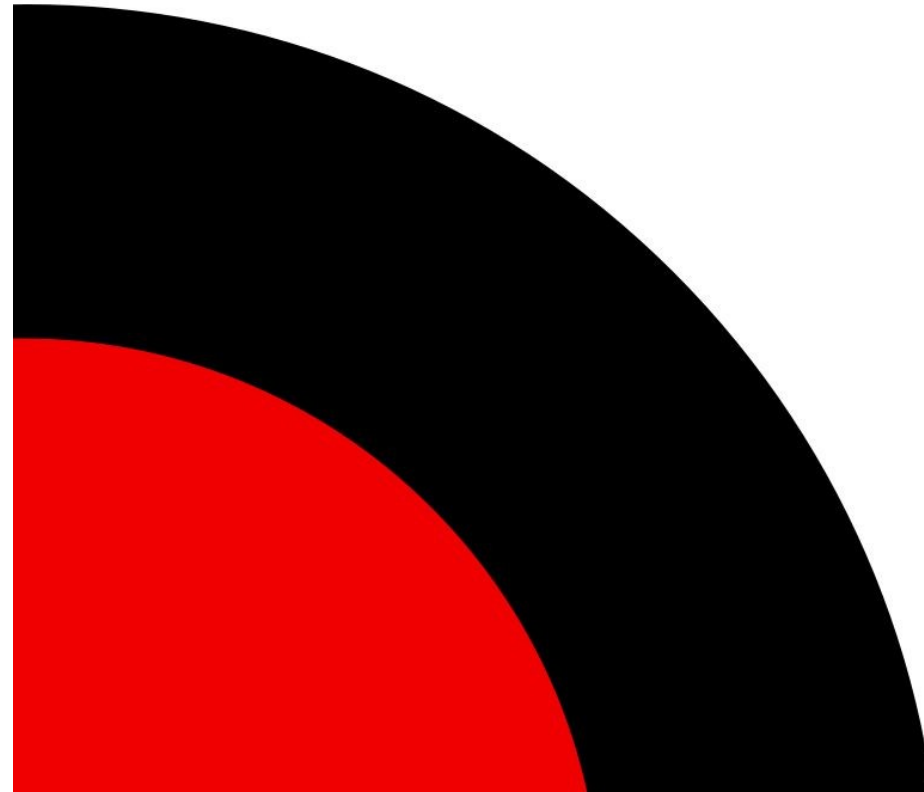
Until one day, the boundary gives way.

And, that dent you've made is called a Ph.D.



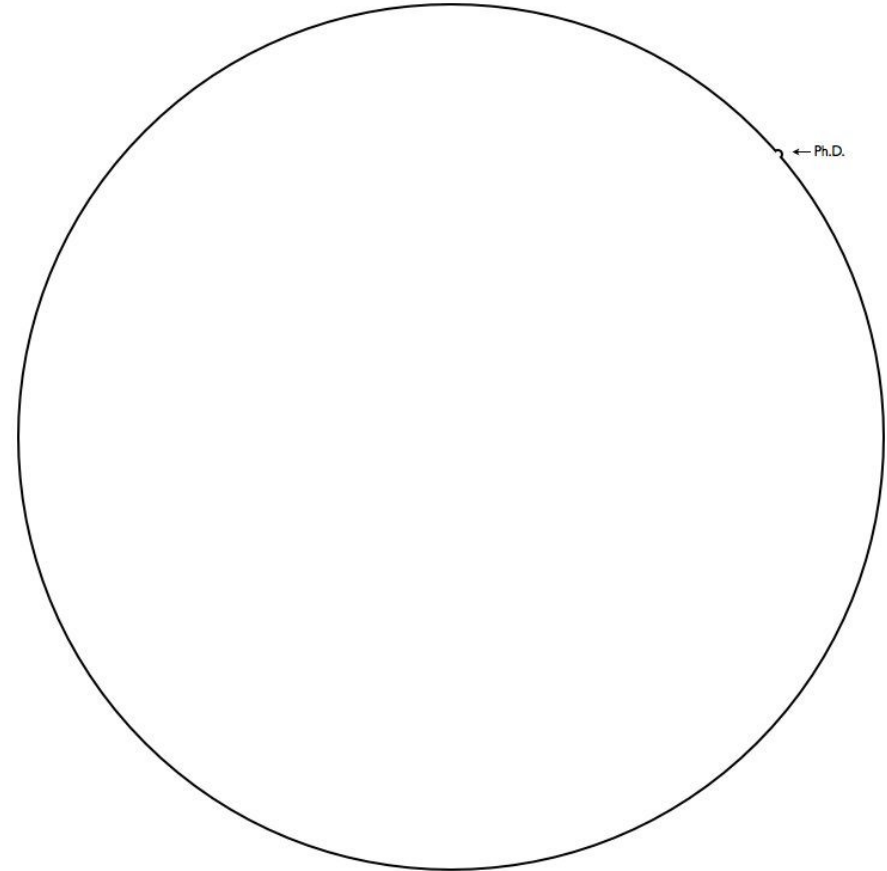
Human knowledge

Of course, the world looks different to you now.



Human knowledge

So, don't forget the bigger picture.



- A PhD requires
 - A PhD student
 - One (or more) PhD advisor
 - A research topic
 - A PhD grant

- The student and advisor should get along well

- The student should be interested about the research topic

- Master projects/internships are a good way to verify it

Getting a PhD grant

- PhD grant : 120k€ over 3 years
 - Includes a salary (~2000€ per month) + social contributions
- Types of PhD grants
 - Institution funding (AMX, Hi!Paris, IPParis, hosting lab, ...)
 - Apply in ~april, results in may/june
 - Research project (ANR, Horizon Europe, ...)
 - Apply in october, results in july
 - CIFRE funding
 - PhD in the industry

- A thesis = a set of research contributions
 - Usually summarized in publications
- What is a publication ?
 - Written article (usually around 10 pages) that treats one particular problem
 - Peer-reviewed
 - Published in a venue (journal, conference, workshop)

Scaling Distributed Deep Learning Workloads beyond the Memory Capacity with KARMA

Mohamed Wahib^{*§}, Haoyu Zhang[†], Truong Thao Nguyen^{*}, Aleksandr Drozd[§], Jens Domke[§], Lingqi Zhang[‡], Ryousei Takano^{*}, Satoshi Matsuoka^{§‡}

^{*} National Institute of Advanced Industrial Science and Technology, Japan

{mohamed.attia,nguyen.truong,takano-ryousei}@aist.go.jp

[†] miHoYo Inc. (This work was done while the coauthor worked in Tokyo Institute of Technology) lynkzhang@gmail.com

[‡] Tokyo Institute of Technology, Tokyo, Japan zhang.l.ai@m.titech.ac.jp

[§] RIKEN Center for Computational Science, Kobe, Japan {aleksandr.drozd,jens.domke}@riken.jp,matsu@acm.org

Abstract—The dedicated memory of hardware accelerators can be insufficient to store all weights and/or intermediate states of large deep learning models. Although model parallelism is a viable approach to reduce the memory pressure issue, significant modification of the source code and considerations for algorithms are required. An alternative solution is to use out-of-core methods instead of, or in addition to, data parallelism.

We propose a performance model based on the concurrency analysis of out-of-core training behavior, and derive a strategy that combines layer swapping and redundant recomputing. We achieve an average of 1.52x speedup in six different models over the state-of-the-art out-of-core methods. We also introduce the first method to solve the challenging problem of out-of-core multi-node training by carefully pipelining gradient exchanges and performing the parameter updates on the host. Our data parallel out-of-core solution can outperform complex hybrid model parallelism in training large models, e.g. Megatron-LM and Turning-NLG.

Index Terms—Deep Neural Networks, Out-of-core, GPUs

I. INTRODUCTION

Training Deep Neural Networks (DNNs) is increasingly becoming one of the main HPC workloads. As model and dataset sizes for Deep Learning (DL) become increasingly large, the memory requirement for training Neural Networks (NNs) increases dramatically. Even though the latest generation of Nvidia GPUs have up to 32 GiB (V100), this capacity remains a major bottleneck in a lot of the cases [1]. For example, with a large network such as ResNet-200 [2], the local batch-size for training cannot be larger than six ImageNet samples, and in ResNet-1001 the local batch size drops down to two samples. This problem is also a challenge for models that require tens of billions of parameters [3], [4], at which the model will not fit into a single GPU and programmers are forced to employ complex model partitioning methods [1].

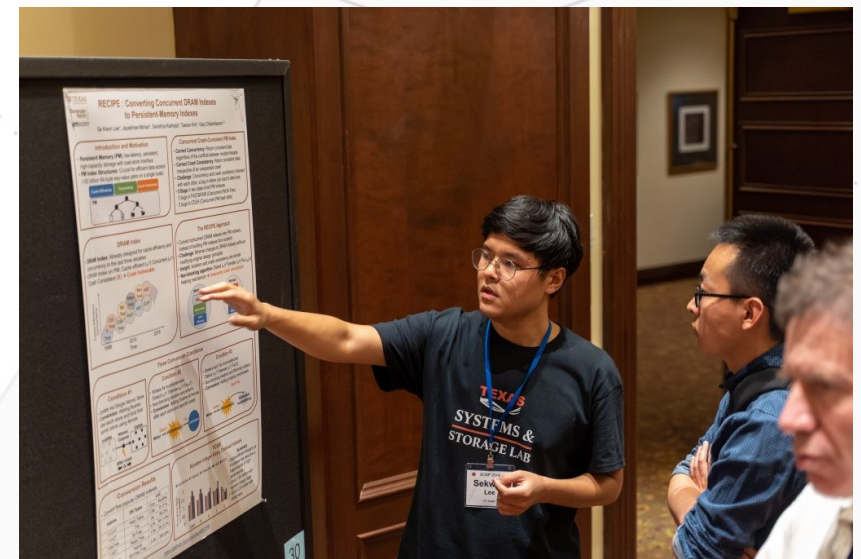
Distributed training can be used if multiple GPUs are available. Data parallelism is a commonly used scheme in which the model is replicated and the training data is distributed. However, this scheme does not reduce memory pressure from model parameters and activations, forcing users to resort to

relatively small, there are cases where even a single training sample is too large to be processed on a single GPU. Such cases include high resolution medical or satellite images which can go up to 2 GiB per sample [5]. In comparison, the widely used ImageNet dataset [6] has images that are smaller than 100 KiB per sample (re-sized to 224×224).

Although model parallelism could be a solution, construction of a cost model and significant modification of the code is needed for every model/dataset/system combination [1], [7], [8]. Another general solution to this memory capacity problem, that we discuss in this paper, is to use out-of-core methods, without or with redundant recompute, to break the GPU memory limitation [9]–[14].

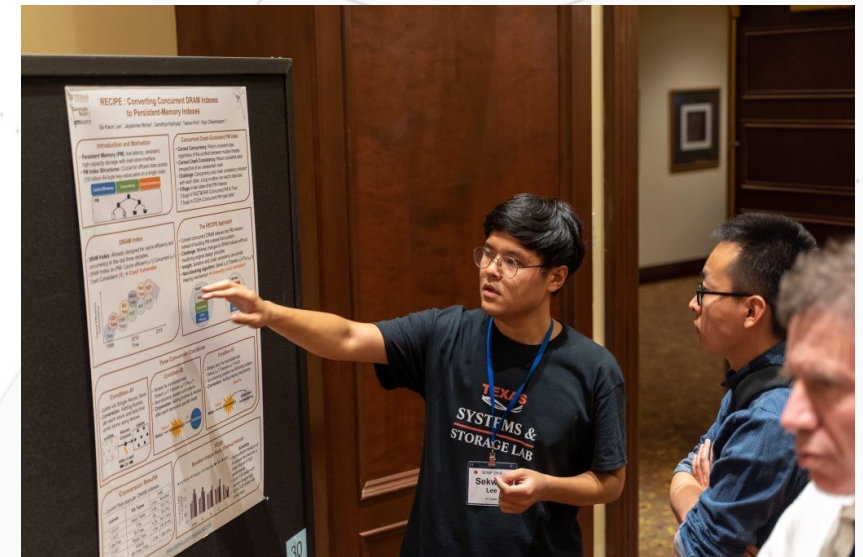
The first challenge that KARMA must address is how to first derive an efficient out-of-core strategy that reduces the stall in the execution pipeline, i.e. address the device occupancy bottleneck. Prefetching and data swapping in general, is a well-researched area. That being said, general prefetching and swapping techniques are not efficient for out-of-core DL since they don't provide a comprehensive considerations of the concurrency requirements and occupancy w.r.t. DL workloads [9], [10]. The main challenge in deriving an efficient prefetching and swapping strategy is to build a robust model for projecting the minimum required concurrency to keep device utilization as close as possible to maximum. This requires taking into consideration specific features and requirements in DL training: reuse of intermediate results from the forward phase in the backward phase, orchestrating complex pipelines in case of distributed training, non-linear dependency between layers, and memory footprint to compute imbalance (i.e. compute is not linearly correlated to the memory footprint). We propose a performance model to derive a capacity-based out-of-core strategy by the means of assuring a minimum concurrency, i.e., available parallelism, that keeps the device at the highest possible utilization. In addition, we identify and utilize any opportunities at which redundantly recomputing layers reduces the stalls in the prefetching pipeline. More specifically, we

- Institution in charge of
 - Assessing the paper quality
 - Publishing the paper
- Main types of venues:
 - Conferences (main publishing venue in CS)
 - International Parallel & Distributed Processing Symposium (IPDPS)
 - Symposium on Operating Systems Principles (SOSP)
 - Workshops (small conference on a specific topic)
 - International Workshop on OpenMP (IWOMP)
 - International Workshop on Runtime and Operating Systems for Supercomputers (ROSP)
 - Demo/Poster session in conferences
 - Journals



Assessing a venue « quality »

- All venues exist in different « qualities »
 - It is very hard to get a paper into a good venue
- Identifying good venues
 - Your advisor knows
 - Look at your list of references
 - Conference/journal ranking
 - CORE (A*, A, B, C) <http://portal.core.edu.au/conf-ranks/>
 - Google Scholar (based on H-index) https://scholar.google.com/citations?view_op=top_venues&vq=eng



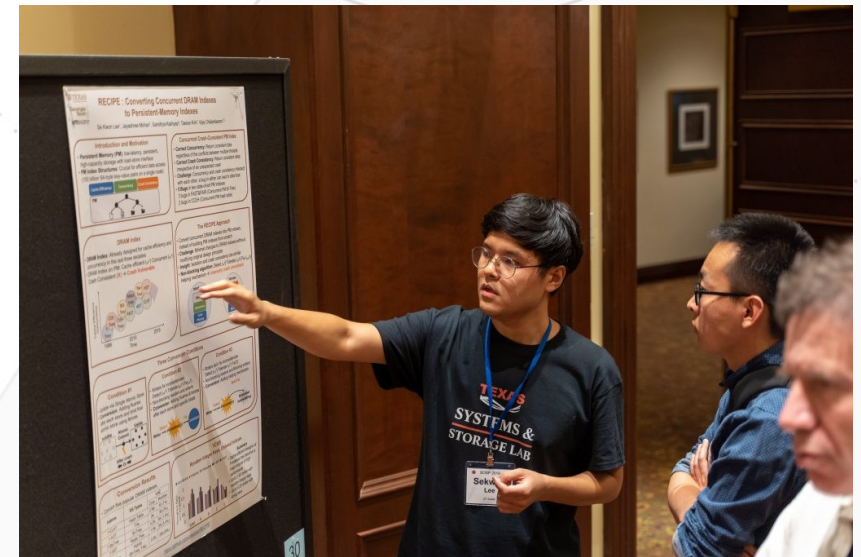
(some of the) top conferences related to PDS / HPDA

- System : SOSP, OSDI, ASPLOS, EuroSys, Usenix ATC, DSN, VEE
- Distributed systems: PODC, ICDCS, DISC, OPODIS, SSS
- Parallel programming: ISCA, IPDPS, PPOPP, SC, ICPP, EuroPar, ICPADS, PDP
- AI: KDD, NeurIPS, AAAI, MLsys

Note on the rank of venues

- **Rank of a venue != importance of a paper**
- A paper published in a low rank venue may be good/important
 - *StarPU: A Unified Platform for Task Scheduling on Heterogeneous Multicore Architectures*
 - Published at EuroPar 2009 (rank B conference)
 - 1896 citations, at the base of task-based parallelism
 - *hwloc: A generic framework for managing hardware affinities in HPC applications*
 - Published at PDP 2010 (rank C conference)
 - 591 citations, used in most HPC applications / runtime systems
- A paper published in a good venue may be insignificant
 - *NewMadeleine: An efficient support for high-performance networks in MPICH2 -- Trahay et al.*
 - Published at IPDPS 2009 (rank A conference)
 - 9 citations

- Write a paper in a specific format
 - Submit the paper for peer review
 - Conferences/workshops: submit before a deadline
 - Journals: submit anytime
 - Paper is reviewed by anonymous experts
 - Assess paper weakness/strength
 - Grade the paper (accept/reject)
- decision based on several reviews (usually 3+)
- Conferences: accept (~20%) / reject (~80%)
 - Journals: reject / major revision / minor revision / accept

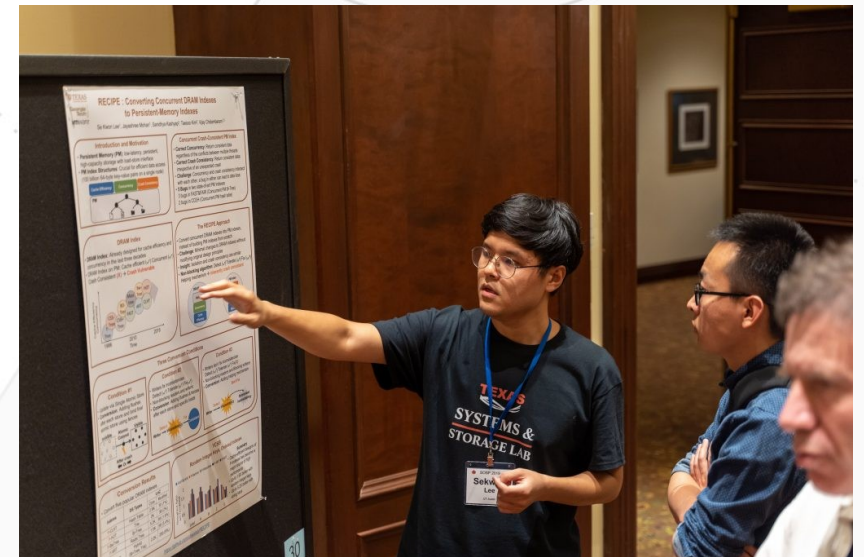


Once the paper is accepted

- The paper is published
 - In a journal
 - As part of the conference proceedings

- You present your work at a conference
 - Travel to a ± fancy place
 - 1 week conference + workshop dedicated to a research topic
 - ~20 minutes presentation of your work

- Your paper gets used/compared to by other researchers



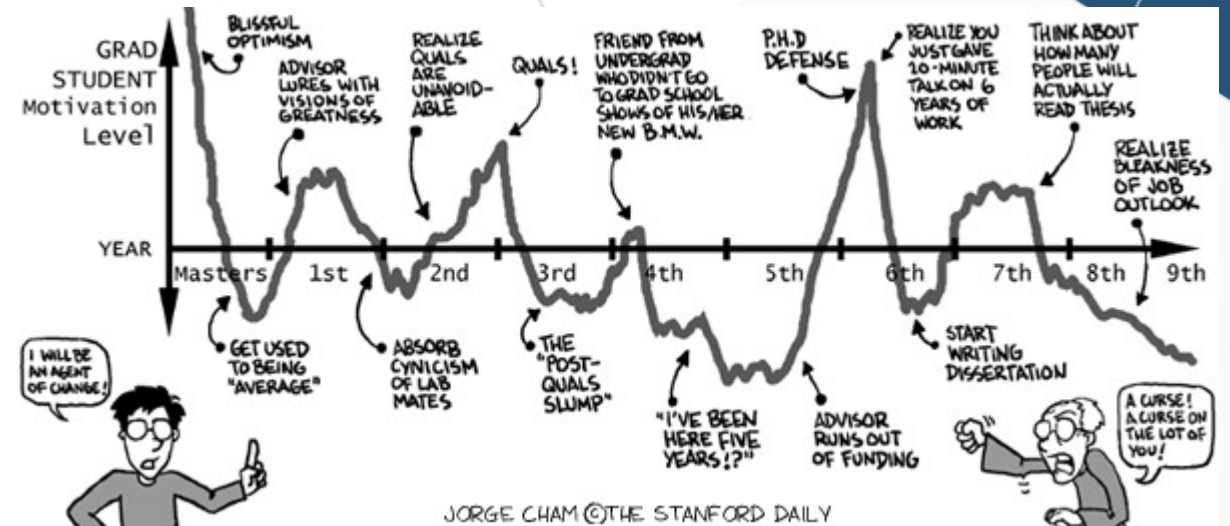
Frustration

- PhD work can be very frustrating and psychologically demanding

- Rejected papers
- Comparison with others can be depressing (Impostor syndrome)
- Stalled development

- Seek help

- Talk with you advisor
- Talk with your colleagues / friend / family
- Universities employ trained psychologist to help you
 - Sylvie Coussot <sylvie.coussot@ip-paris.fr >
- 30 % of PhD students seek help¹



- Typical career in academia in France

- PhD: 3.5 years – ~2000 € / month
- Post-doc abroad: 2 years – 2500+ € / month
- Associate professor (*maître de conférence*) / researcher (*chargé de recherche*) : 10 years – 2000+ € / month
- Full professor (*professeur des universités*) / senior researcher (*directeur de recherche*) : 25 years – 3000+ € / month

- After a PhD in Computer Science¹

- In France : 73% / overseas : 27 %
- Permanent position : 73% / fixed-term contract : 27 %
- Academia : 28% / Industry : 55% / Other : 17 %
- Median salary (in France): 41538 € / 43077 € (industry)

¹ https://enquetes-sphinx.u-psud.fr/PHdFutur/PhDs_Future/PhDs_Future.htm
Census among 233 PhD obtained in 2016-2018 from Université Paris-Saclay