

Michele (pirroh) Catasta

ABOUT ME

Head of Applied Research at Google Labs

working on *AI applied to Source Code, Large Language Models*

Personal page: <https://pirroh.fyi>

LinkedIn profile: <https://linkedin.com/in/pirroh>

Twitter profile: <https://twitter.com/pirroh>

GitHub profile: <https://github.com/pirroh>

Google Scholar: <https://scholar.google.com/citations?user=Zxf0SWMAAAAJ&hl=en&oi=ao>

EMAIL

✉ pirroh@cs.stanford.edu (personal)

✉ pirroh@google.com (work)

AREAS OF EXPERTISE

AI4Code, Machine Learning, Data Science, Information Retrieval

EDUCATION

STANFORD UNIVERSITY

Postdoc and Instructor in Machine Learning

2017 – 2019

- Advisor: Prof. Jure Leskovec
- Affiliated with *Statistical Machine Learning Group*, SNAP and InfoLab

EPFL, SWITZERLAND

Ph.D. in Computer Science

2010 – 2015

- Advisors: Prof. Karl Aberer and Prof. Philippe Cudre-Mauroux
- Thesis: *MEMORIES: Memory-based Information Systems*

NATIONAL UNIVERSITY OF IRELAND, GALWAY

M.Sc. with Research Honors in Applied Computing and IT

2007 – 2009

- Advisors: Prof. Stefan Decker and Dr. Giovanni Tummarello
- Thesis: *Scalability Solutions for the Web of Data*
- Full tuition scholarship and Research Assistantship contract

EXPERIENCE

IRREGULAR EXPRESSIONS

Limited Partner and Advisor

2022 – Present

IrregEx is a group of seasoned SWE and product leaders advising and investing in early-stage startups focused on developer tools.

X, THE MOONSHOT FACTORY [FORMERLY GOOGLE X] (MOUNTAIN VIEW, CA, USA)

Head of Applied Research

2021 – 2022

AI applied to Source Code, Large Language Models

STANFORD UNIVERSITY (STANFORD, CA, USA)

Research Scientist and Instructor

2020

Contributed to several projects: Graph Neural Networks and Transformer architectures applied to relational data (with a focus on Source Code), Open Graph Benchmark, recommender system for Wikipedia editors, resource allocation on crowdsourcing platforms (IARPA HFC), machine common sense (DARPA MCS).

EPFL (LAUSANNE, SWITZERLAND)

Scientific Collaborator and Lecturer

2015 – 2017

Created and taught the first edition of Applied Data Analysis (<https://ada.epfl.ch>), a brand new course in the EPFL Data Science specialization. Managed several research projects with a focus on machine learning, recommender systems, data mining, etc.

MIT MEDIA LAB (CAMBRIDGE, MA, USA)

Visiting PhD in Prof. Sandy Pentland's group

Summer 2014

Implemented a gesture-based authentication prototype on iOS, used by the Human Dynamics lab.

GOOGLE (MOUNTAIN VIEW, CA, USA)

Software Engineer Intern, Infrastructure Storage

Summer 2011

Given the resource-consumption information and the SLAs of the storage layer, I worked on a compiler that translates high-level description of SLOs into low-level monitoring code (using Python) and on prediction models for timeseries (using Dremel).

YAHOO RESEARCH (BARCELONA, SPAIN)

Student Researcher in Prof. Ricardo Baeza-Yates' group

Summer 2008

Developed a Big Data processing pipeline (using Hadoop and Pig) to extract metadata (such as `schema.org` tags) from the Yahoo Web crawl.

DERI (GALWAY, IRELAND)

Research Assistant

2007 – 2009

Co-founder of Sindice.com, the largest Semantic Web search engine at the time. I designed and developed several key components of the backend, with a focus on the scalability challenges (using MapReduce, Hadoop, Lucene, etc.) The core technologies developed for Sindice.com evolved into an investigative intelligence platform (<https://siren.io>) which—to date—has secured \$15M+ in funding.

GRANTS AND
AWARDS

- ◇ Granted \$5,000 in Google Cloud Platform credits for the research project ACADEMIC DIASPORA, 2020
- ◇ Nominated **Google Cloud Faculty Expert**, 2020
- ◇ Received an Education Grant for \$195,000 in Google Cloud Platform credits to support the course CS224W: MACHINE LEARNING WITH GRAPHS, 2019
- ◇ Granted \$15,000 in Amazon Web Services credits for the research project WEBGENOME, 2018
- ◇ Visiting Scholar grant from the Information School of the University of Sheffield, 2016
- ◇ Google Student grant for the World Wide Web conference, Firenze, Italy, 2015
- ◇ Short Visit grant for “Evaluating Information Access Systems”, University of Sheffield, 2015
- ◇ Invited to the Global Young Scientists Summit @one-north, Singapore, 2015
- ◇ Lead author of the CROSS grant (60,000 CHF), multidisciplinary collab. with UNIL, 2014
- ◇ 3rd prize at the Semantic Web Challenge, ISWC, 2013
- ◇ Lead author of the MemorySense grant (\$150,000), collab. with Samsung Research USA, 2013
- ◇ Outstanding Teaching Assistant Award, EPFL, 2012
- ◇ Hypothes.is Reputation Fellow, 2012
- ◇ 3rd prize at the Semantic Web Challenge, ISWC, 2009
- ◇ 2nd prize at the Asian Semantic Web Conference, Industrial Track, 2009
- ◇ 1st prize at the European Semantic Web Technology Conference, Business Idea Contest, 2008
- ◇ Saltlux Best Paper Award for the Sindice journal paper, 2007

TEACHING
ACTIVITIES

- ◇ **Advisor** for CS329S: MACHINE LEARNING SYSTEMS DESIGN at Stanford University, 2021
- ◇ **Co-instructor** for CS224W: MACHINE LEARNING WITH GRAPHS (≈ 200 students) and CS246: MINING MASSIVE DATA SETS (≈ 300 students) at Stanford University, 2018 – 2020
- ◇ Mentor at LEAD THE FUTURE, a network to support Italian excellence in STEM, 2020 – Present
- ◇ **Instructor** for CS341: PROJECT IN MINING MASSIVE DATA SETS at Stanford University, 2019
- ◇ Creator and first instructor of APPLIED DATA ANALYSIS (≈ 250 students), a brand new course in Data Science at EPFL, **second largest course offered by the CS department**, 2016 – 2017
- ◇ Core member of the committee which developed the Data Science MSc at EPFL, 2015 – 2017
- ◇ Lecturer for the course DISTRIBUTED INFORMATION SYSTEMS at EPFL, 2015 – 2016
- ◇ Lead TA for the course DISTRIBUTED INFORMATION SYSTEMS at EPFL, 2011 – 2014
- ◇ TA for the course INTRODUCTION TO DATABASE SYSTEMS at EPFL, 2012
- ◇ Co-founder of “HACKERS AT EPFL”, organizing hackathons, tech talks, etc., 2012 – 2015

ACADEMIC
ACTIVITIES

- ◇ Co-organizer of the MINING AND LEARNING WITH GRAPHS workshop at KDD, 2017
- ◇ Co-organizer of the Dagstuhl Seminar CROWDSOURCING RESEARCH - TRANSCENDING DISCIPLINARY BOUNDARIES, 2016

- ◇ **Senior PC member** of CIKM 2020, CIKM 2021
- ◇ **PC member** of SocInfo 2012, CIKM 2014, ESWC 2015, ICWE 2016, ESWC 2016, ISWC 2016 Doctoral Consortium, HCOMP 2016, WWW 2017 (and posters), Wiki Workshop 2017, ESWC 2017, KDD 2017, WebSci 2017, ISWC 2017 Doctoral Consortium, CIKM 2017, AAAI 2018, WSDM 2018, WWW 2018 (and posters), Wiki Workshop 2018, WebSci 2018, KDD 2018 Applied Data Science, CIKM 2018, ICDM 2018, HumL at WWW 2018 and ISWC 2018, CrowdBias Workshop at HCOMP 2018, Euro CSS Symposium 2018, WSDM 2019, Wiki Workshop 2019, KDD 2019, WWW 2020, ECML-PKDD 2020, GRL+ Workshop @ ICML 2020, ICML 2020, NeurIPS 2020, ICLR 2021, AAAI 2021, ICML 2021, NeurIPS 2021, ICLR 2022, ICML 2022.
- ◇ **Reviewer** for ACM Transactions on Social Computing, IEEE Transactions on Knowledge and Data Engineering, EPJ Data Science, Nature Scientific Data, Harvard Data Science Review.

INVITED TALKS

- ◇ *AI meets Source Code: status quo and outlooks*, SYNAPSE AI SYMPOSIUM, 2022
- ◇ *AI meets Source Code: status quo and outlooks*, EPFL, hosted by Prof. Robert West, 2022
- ◇ *AI meets Source Code: status quo and outlooks*, UNIVERSITY OF FRIBOURG, hosted by Prof. Philippe Cudre-Mauroux, 2022
- ◇ *Moderator of Panel: Big Data and AI in the Classroom with Google Cloud*, GOOGLE FACULTY INSTITUTE, 2020
- ◇ *Advancements in Graph Deep Learning*, GRAPH THE PLANET, Bloomberg SF, 2020
- ◇ *Structuring Wikipedia Articles with Section Recommendations*, UNIVERSITY OF FRIBOURG, hosted by Prof. Philippe Cudre-Mauroux, 2018
- ◇ *Recommender Systems in the Wild*, PAPERS WE LOVE, Square SF, 2017
- ◇ *Big Data Mining*, TUI STRATEGY WORKSHOP IN AI AND ML, Palma de Mallorca, 2017
- ◇ *Recommender Systems in the Wild*, UNIVERSITÀ SVIZZERA ITALIANA, hosted by Prof. Crestani, 2017
- ◇ *Lecturer at the CUSO Winter School “Data Science in Information Society: From Data Acquisition to Data Analysis”*, 2017
- ◇ *Entity Typing on the Web*, SAPIENZA UNIVERSITY OF ROME, hosted by Prof. Roberto Navigli, 2016
- ◇ *MEM0R1ES: Memory-based Information Systems*, UNIVERSITY OF LEEDS, hosted by Prof. Ruddle, 2016
- ◇ *MEM0R1ES: Memory-based Information Systems*, UNIVERSITY OF SHEFFIELD, hosted by Prof. Gianluca Demartini, 2016
- ◇ *Entity Typing on the Web*, YAHOO LABS LONDON, hosted by Dr. Peter Mika, 2015
- ◇ *Invited lecture for the Technology Ventures in IC course at EPFL*, hosted by Prof. Bugnion, 2015
- ◇ *MEmoIt: don’t write your diary, sense it*, DBTA Workshop on “Life Logging and Long-Term Digital Preservation” at UNIVERSITY OF LUGANO, 2015
- ◇ *MemorySense: let your smartphone figure out the memorable moments of your life*, SAMSUNG RESEARCH AMERICA, 2014

SELECTED PUBLICATIONS

- A Chowdhery, S Narang, J Devlin, M Bosma, G Mishra et al. “PALM: SCALING LANGUAGE MODELING WITH PATHWAYS”, *arXiv:2204.02311*
- D Zugner, T Kirschstein, M Catasta, J Leskovec, S Gunnemann, “LANGUAGE-AGNOSTIC REPRESENTATION LEARNING OF SOURCE CODE FROM STRUCTURE AND CONTEXT”, *ICLR 2021*
- T Vaucher, A Spitz, M Catasta, R West, “A DECADE IN QUOTES: MINIMALLY SUPERVISED QUOTATION ATTRIBUTION IN MASSIVE NEWS CORPORA”, *WSDM 2021*, **Oral**
- W Hu, M Fey, M Zitnik, Y Dong, H Ren, B Liu, M Catasta, J Leskovec “OPEN GRAPH BENCHMARK: DATASETS FOR MACHINE LEARNING ON GRAPHS”, *NeurIPS 2020*, **Spotlight**
- T Piccardi, M Catasta, L Zia, R West, “STRUCTURING WIKIPEDIA ARTICLES WITH SECTION RECOMMENDATIONS”, *SIGIR 2018*
- J Rappaz, M Catasta, R West, K Aberer, “LATENT STRUCTURE IN COLLABORATION: THE CASE OF REDDIT R/PLACE”, *ICWSM 2018*
- J Rappaz, ML Vladarean, J McAuley, M Catasta, “BARTERING BOOKS TO BEERS: A RECOMMENDER SYSTEM FOR EXCHANGE PLATFORMS”, *WSDM 2017*
- JE Ranvier, M Catasta, M Vasirani, K Aberer, “ROUTINESense: A MOBILE SENSING FRAMEWORK FOR THE RECONSTRUCTION OF USER ROUTINES”, *Mobiquitous 2015*, **Best Paper Award**

DE Difallah, M Catasta, G Demartini, P Ipeirotis, P Cudre-Mauroux, “THE DYNAMICS OF MICRO-TASK CROWD-SOURCING – THE CASE OF AMAZON MTURK”, *WWW 2015*

DE Difallah, M Catasta, G Demartini, P Cudre-Mauroux, “SCALING-UP THE CROWD: MICRO-TASK PRICING SCHEMES FOR WORKER RETENTION AND LATENCY IMPROVEMENT”, *HCOMP 2014*

M Catasta, A Tonon, DE Difallah, G Demartini, K Aberer, P Cudre-Mauroux, “TRANSACTIVEDB: TAPPING INTO COLLECTIVE HUMAN MEMORIES”, *VLDB 2014*

K Aberer, M Catasta, G Christodoulou, I Gavrilovic, F Hrisafov, M Monney, A Ouazaki, B Perovic, H Radu, JE Ranvier, M Vasirani, Z Yan, “MEMO-IT: DON’T WRITE YOUR DIARY, SENSE IT”, *UbiComp 2014*

M Catasta, A Tonon, DE Difallah, G Demartini, K Aberer, P Cudre-Mauroux, “HIPPOCAMPUS: ANSWERING MEMORY QUERIES USING TRANSACTIVE SEARCH”, *WWW 2014*

K Aberer, M Catasta, H Radu, JE Ranvier, M Vasirani, Z Yan, “MEMORYSENSE: RECONSTRUCTING AND RANKING USER MEMORIES ON MOBILE DEVICES”, *PerCom 2014*

A Tonon, M Catasta, G Demartini, P Cudre-Mauroux, K Aberer, “TRANK: RANKING ENTITY TYPES USING THE WEB OF DATA”, *ISWC 2013*, **Best Paper Award nominee**

LH Tran, M Catasta, LK McDowell, K Aberer, “NEXT PLACE PREDICTION USING MOBILE DATA”, *Mobile Data Challenge (by Nokia) Workshop 2012*

R Delbru, N Toupikov, M Catasta, G Tummarello, “A NODE INDEXING SCHEME FOR WEB ENTITY RETRIEVAL”, *ESWC 2010*

R Delbru, N Toupikov, M Catasta, G Tummarello, S Decker, “HIERARCHICAL LINK ANALYSIS FOR RANKING WEB DATA”, *ESWC 2010*

JOURNAL PAPERS T Sinha, M Kapur, R West, M Catasta, M Hauswirth, D Trninic “DIFFERENTIAL BENEFITS OF EXPLICIT FAILURE-DRIVEN AND SUCCESS-DRIVEN SCAFFOLDING IN PROBLEM-SOLVING PRIOR TO INSTRUCTION”, *Journal of Educational Psychology*, 2020

A Tonon, M Catasta, R Prokofyev, G Demartini, K Aberer, P Cudre-Mauroux “CONTEXTUALIZED RANKING OF ENTITY TYPES BASED ON KNOWLEDGE GRAPHS”, *Journal of Web Semantics*, 2016

M Catasta, A Tonon, V Pasquier, G Demartini, P Cudre-Mauroux, K Aberer, “B-HIST: ENTITY-CENTRIC SEARCH OVER PERSONAL WEB BROWSING HISTORY”, *Journal of Web Semantics*, 2014

G Tummarello, R Cyganiak, M Catasta, S Danielczyk, R Delbru, S Decker, “SIG.MA: LIVE VIEWS ON THE WEB OF DATA”, *Journal of Web Semantics*, 2010, **Top-25 Most Cited article of JoWS**

E Oren, R Delbru, M Catasta, R Cyganiak, H Stenzhorn, G Tummarello, “SINDICE.COM: A DOCUMENT-ORIENTED LOOKUP INDEX FOR OPEN LINKED DATA”, *International Journal of Metadata, Semantic and Ontologies*, 2008

BOOK CHAPTERS JE Ranvier, M Catasta, I Gavrilovic, G Christodoulou, H Radu, K Aberer, “MEMOIT: FROM LIFELOGGING APPLICATION TO RESEARCH PLATFORM”, Book chapter in *Mobile Application Development, Usability, and Security*, IGI Global, 2016

R Delbru, N Toupikov, M Catasta, R Fuller, G Tummarello, “SIREN: EFFICIENT SEARCH ON SEMI-STRUCTURED DOCUMENTS”, Book chapter in *Lucene in Action 2nd Edition* (In Action series), Manning Publications Co., 2009

M Catasta, R Delbru, N Toupikov, G Tummarello, “MANAGING TERABYTES OF WEB SEMANTICS DATA”, Book chapter in *Semantic Web Information Management*, Springer, 2009

BOOK G Demartini, DE Difallah, U Gadiraju, M Catasta, “AN INTRODUCTION TO HYBRID HUMAN-MACHINE INFORMATION SYSTEMS”, Book in *Foundations and Trends in Web Science*, NOW, 2017

TUTORIALS
◇ “It’s Getting Crowded!”: How To Use Crowdsourcing Effectively for Web Science Research, <https://itsgettingcrowded.wordpress.com>, *WebSci 2016*
◇ Using Crowdsourcing Effectively for Social Media Research, *ICWSM 2016*