# chandan singh

✉ cs1@berkeley.edu    ◐ csinva    🎓
◉ csinva.io    in csinva    🐦 csinva

## education

**phd  | machine learning**
uc berkeley | '17-'22
research: interpretable ml
advisor: bin yu
collaborators:
    a. kornblith (medicine)
    s. upadhyayula (biology)

**bs  | cs & math**
university of virginia | '14-'17
double major

## skills

language models • deep learning
data science • data cleaning
huggingface • pytorch
rule-based models • causal inference

## awards

berkeley grad slam semifinalist '19, '22
pdsoros fellowship finalist '19
outstanding teaching award '18
uva rader research award '17
uva undergrad symposium winner '17
raven honor society '16-'17
icpc regional qualification '14-'16
1st place microsoft code jam '16
3rd place google games uva '17
2nd place apt puzzle competition '17
rodman scholarship '14-'17

## teaching

berkeley | summer 2018
machine learning: cs 189/289 ✆
lectures to class of 80+ students

berkeley | fall 2019
artificial intelligence: cs 188 ✆

## service

basis education volunteering '19-'22
bair undergrad mentoring '18-'22
computer literacy volunteering '15-'17
neurips reviewer '20, '23
acl reviewer '22
iclr,cvpr,aaai,neurips reviewer '21

## experience

### microsoft research
senior researcher (deep learning lab) | summer '22 - present
- improving the interpretability of language models
- researching scientific/medical knowledge discovery with language models

### health tech
paige ai | research scientist | summer '21 - summer '22
- interpretable deep learning in digital pathology (especially bladder cancer)

response4life | volunteer data scientist | spring '20
- helped develop, integrate, and deploy models to forecast covid-19 severity

pacmed ai | healthcare ml internship | summer '19
- developed interpretable, tabular machine-learning models for healthcare

### phd
berkeley | interpretable ml research (bin yu lab ✆) | fall '17 - spring '22
- developed post-hoc interpretation methods for ml models (e.g. neural nets)
- developed interpretable models in medicine, biology, and computer vision

aws | ml fairness internship (pietro perona lab ✆) | summer '20
- testing for bias with causal matching using GANs

meta ai | computer vision internship | summer '17
- investigated unsupervised deep learning for segmentation of satellite imagery

### undergrad
hhmi | ml research (srini turaga lab ✆) | summer '14, '15, '16
- researched neural image segmentation and biophysical simulations

uva | ml research (yanjun qi lab ✆) | fall '16 – spring '17
- developed multi-task graphical models for analyzing functional brain connectivity

uva | comp. neuroscience research (william levy lab ✆) | fall '14 - fall '16
- developed biophysical models of single-neuron computation
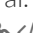
## selected publications

interpretability × language models
- augmenting interpretable models with llms **cs**, et al. *nature comm.*, '23 ✆ </>
- tree prompting morris*, **cs***, rush, gao, & deng *emnlp*, '23 ✆ </>
- explaining data patterns in natural language **cs***, morris*, et al. *emnlp workshop*, '23 ✆

interpretability × deep learning
- adaptive wavelet distillation from dnns: ha, **cs**, et al. *neurips '21* ✆ </>
- aligning dnns by regularizing explanations: rieger, **cs**, et al. *icml '20* ✆ </>
- hierarchical interpretations for dnn predictions: **cs***, murdoch*, & yu, *iclr '19* ✆ </>

interpretabity × rules
- imodels: an interpretability package: **cs***, nasseri*, tan, tang, & yu, *joss '21* ✆ </> ◉
- fast interpretable greedy-tree sums: tan*, **cs***, nasseri*, agarwal* et al. *arxiv '22* ✆ </>
- hierarchical tree shrinkage agarwal*, tan*, ronen, **cs**, & yu *icml '22* ✆ </>