

ModelX: A Novel Transfer Learning Approach Across Heterogeneous Datasets

Arunavo Dey¹, Neil Anthony³, Aakash Dhakal¹, Jayaram Thigarajan², Jae-Seung Yeom²,
Tapasya Patki², Tom Scogland², Tanzima Z. Islam¹

¹Texas State University

²Lawrence Livermore National Laboratory

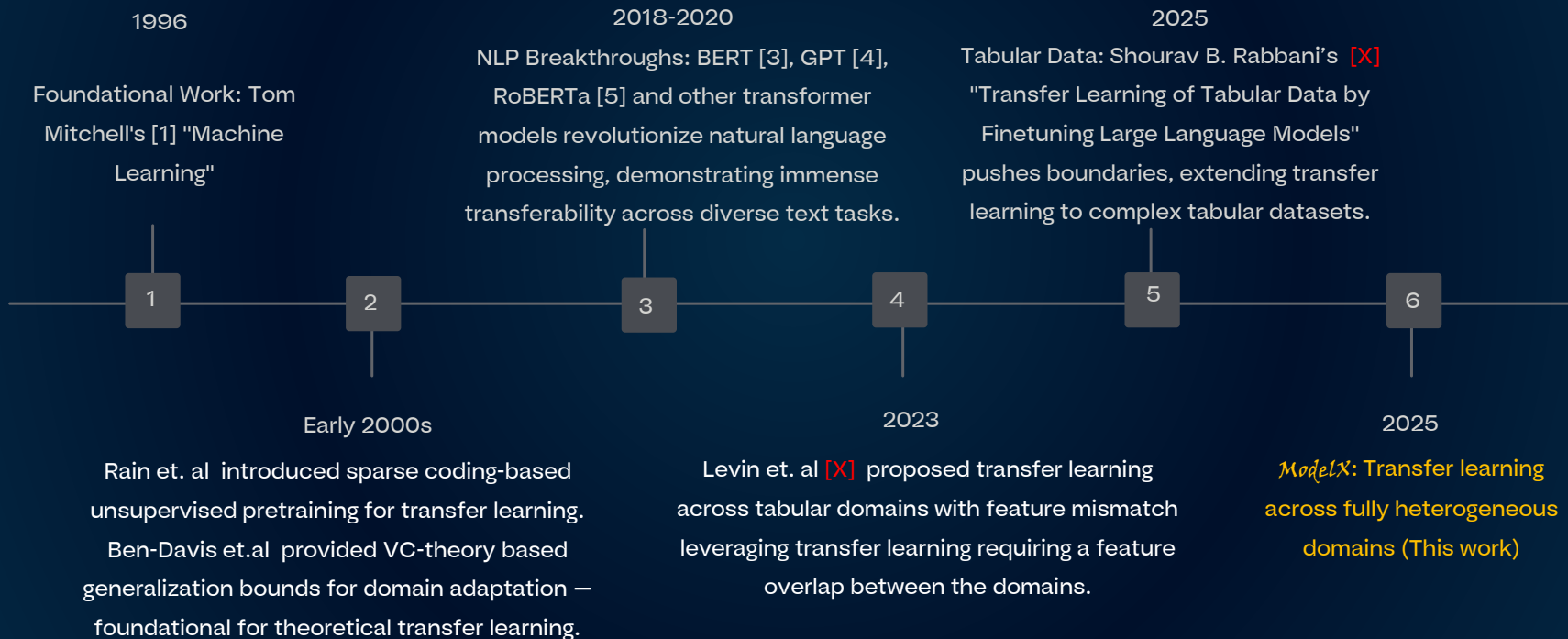
³University of California Santa Barbara



Lawrence Livermore
National Laboratory



The Evolution of Transfer Learning



Where Transfer Learning Excels and Falters

- Similar data distributions across datasets



Similar
Domains

- Need a lot of data for training supervised models



Abundant
Source Data

- Assumes overlap between feature names



Feature
Overlap



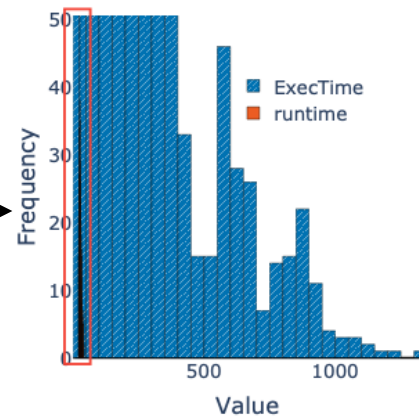
Domain
Divergence



Data Scarcity



Disjoint feature sets



Value



SOTA cannot do it!

Excels

Falters

Where Transfer Learning Excels and Falters

- Similar data distributions across datasets



Similar Domains

- Need a lot of data for training supervised models



Abundant Source Data

- Assumes overlap between feature names



Feature Overlap



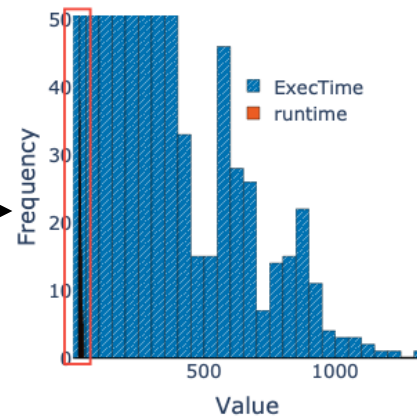
Domain Divergence



Data Scarcity



Disjoint feature sets



Data and structural heterogeneity could be caused by



CoMD on Intel Xeon E5-2695 v2

1	application	algorithm	bw_level	ipath_0	ipath_1	task_count	power_cap	runtime	thread_count
2	CoMD	pak	1	2	0	4096	64	54.4034	24
3	CoMD	pak	1	2	0	4096	80	46.5075	24
4	CoMD	pak	1	2	0	4096	115	42.4216	24
5	CoMD	pak	1	2	0	2048	64	41.7099	16

Kripke on Intel Xeon E5-2695 v2

	DRAMPowerPerNode	ProcessorPowerPerNode	Ranks	Netting	Dset	Gset	CMP	PKG_LIMIT	AvgInst	AvgIpc	AvgArithFpu	AvgFreq	AvgTemp	ProcessorPower	DRAMPower	ExecTime
0	31.572256	86.94462	32	GGZ	8	1	4	50	6.82929e+10	3.7360	6.248262e+08	2.1110	34.2837	1390.4714	505.1561	7.6394
1	32.813519	95.389194	32	GGZ	8	1	4	55	6.088132e+10	4.2230	6.212548e+08	2.3428	33.1925	1526.2271	525.0163	6.0149
2	33.178587	103.280106	32	GGZ	8	1	4	60	5.925558e+10	4.5708	6.208049e+08	2.5125	35.6479	1652.4817	530.8574	5.4116
3	33.871169	113.623712	32	GGZ	8	1	4	65	5.633770e+10	4.8950	6.194195e+08	2.8936	35.8681	1818.0754	541.9387	4.8056

Agenda



Problem and motivation



HPC scenarios & challenges



Methodology



Evaluations



Application of *ModelX* for job scheduling



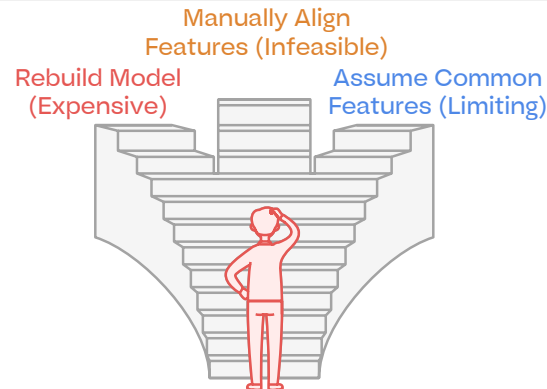
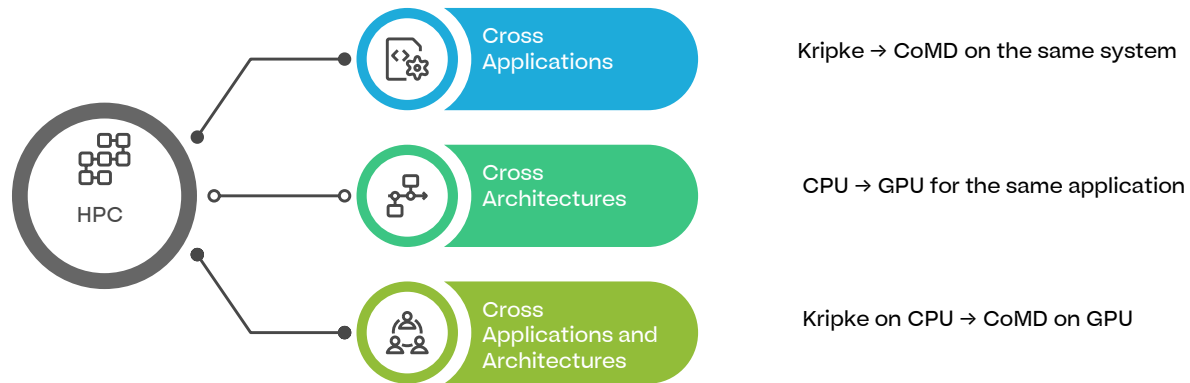
Conclusions & Future work

Transfer Learning Scenarios in HPC are Challenging Due to Heterogeneity

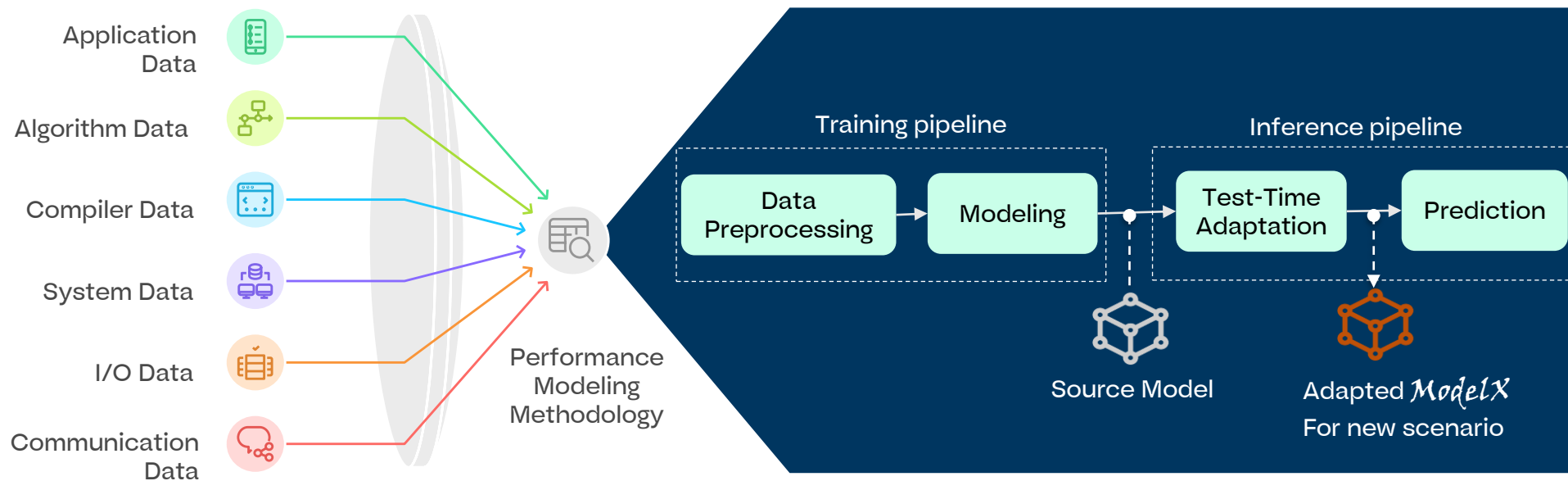
Challenges:

- Domain divergence: Data distribution shift across homogeneous datasets
- Data scarcity: Extensive data collection is expensive
- Feature heterogeneity (different names, counts, order of features)

How do **they** handle domain divergence, feature mismatch or disjoint feature sets today?



Performance Modeling Methodology using ML

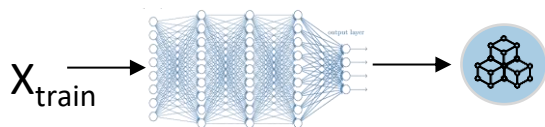


- Assumption: All data sources during **training** are homogeneous.
- No assumption during inference time.

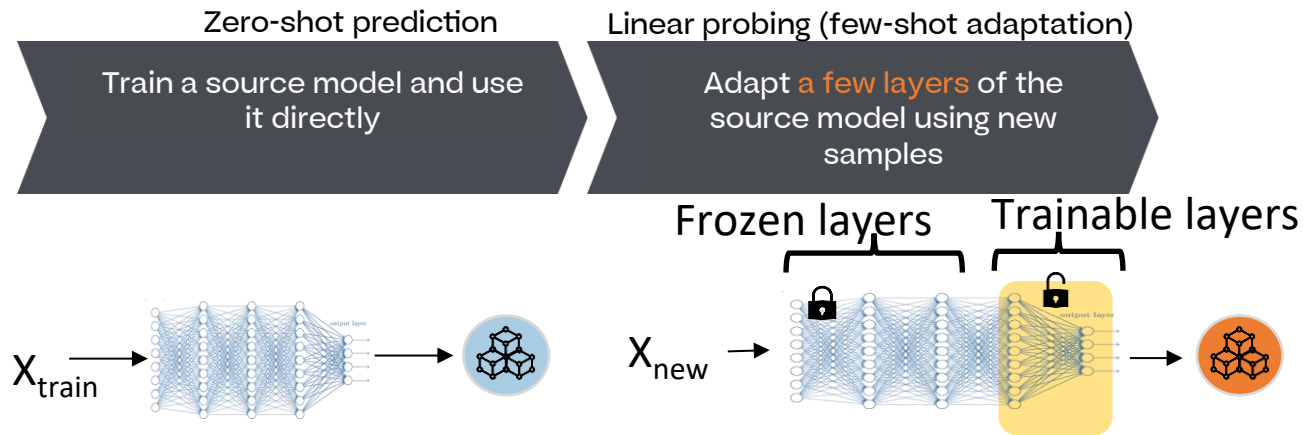
Current Test-Time Adaptation Approaches

Zero-shot prediction

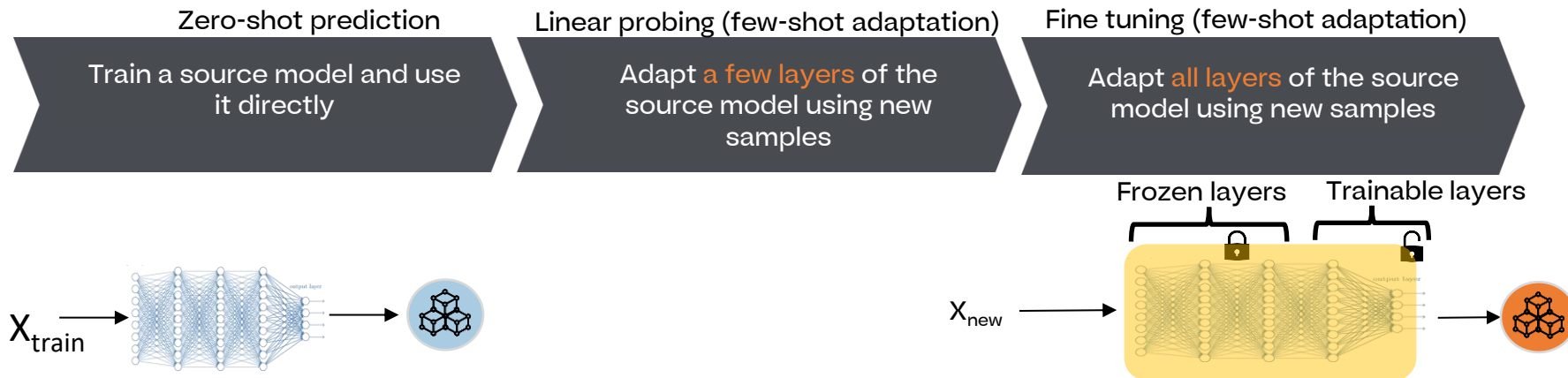
Train a source model and use
it directly



Current Test-Time Adaptation Approaches



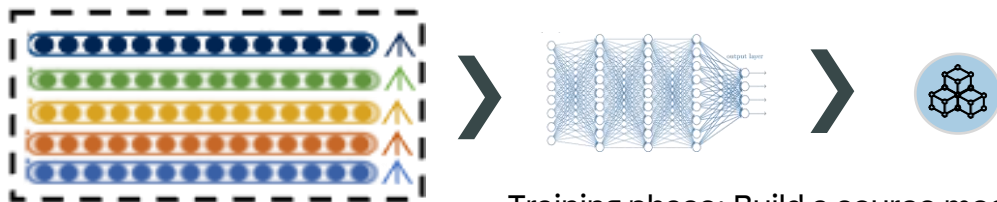
Current Test-Time Adaptation Approaches



Proposed Solution: Bridging Heterogeneity During Inference Time

- Introducing a novel approach for robust performance prediction across diverse domains.

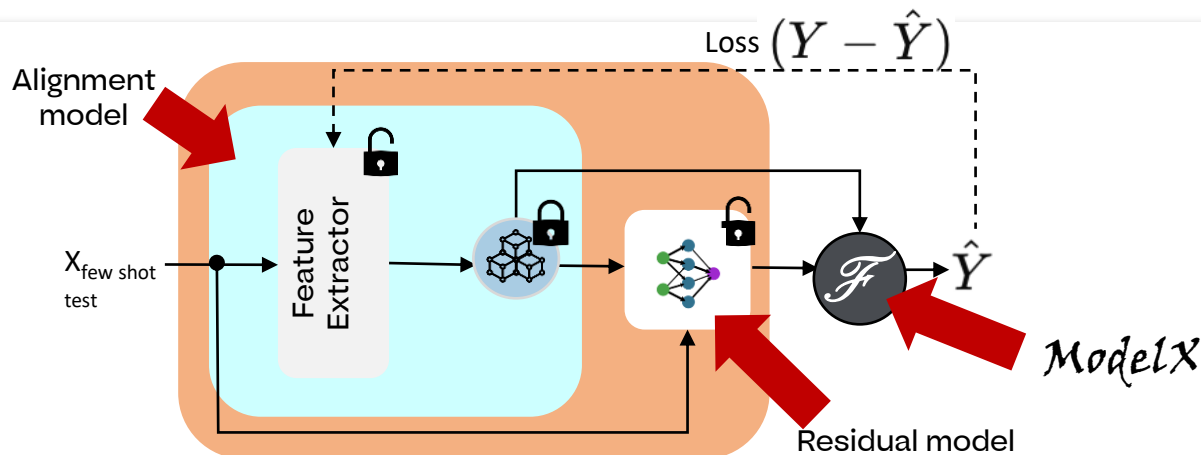
During training: Source model building using homogeneous datasets



Training phase: Build a source model from one or more homogeneous source datasets

During inference:

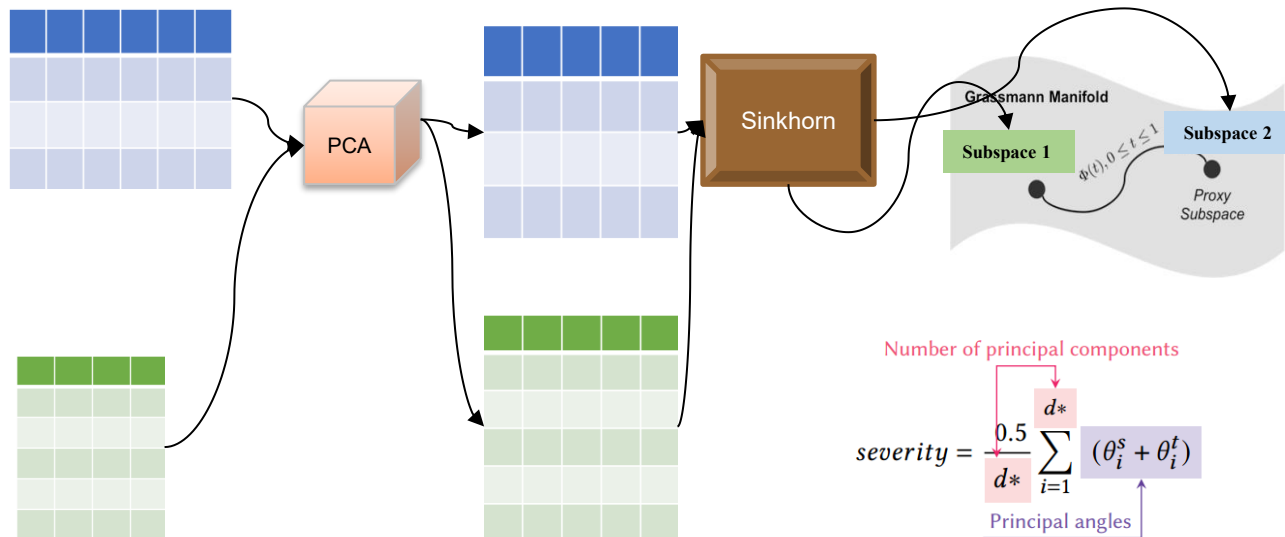
- Step 1: Train a feature extractor network (M_A) using new target few-shot samples using source model's prediction loss
- Step 2: Use just the new target few-shot samples to train an additional residual model M_R



Inference phase: A few-shot test-time adaptation method that can learn from new samples that may have disjoint features, data divergence

Proposed Explainability Measure for Quantifying the Divergence between Datasets

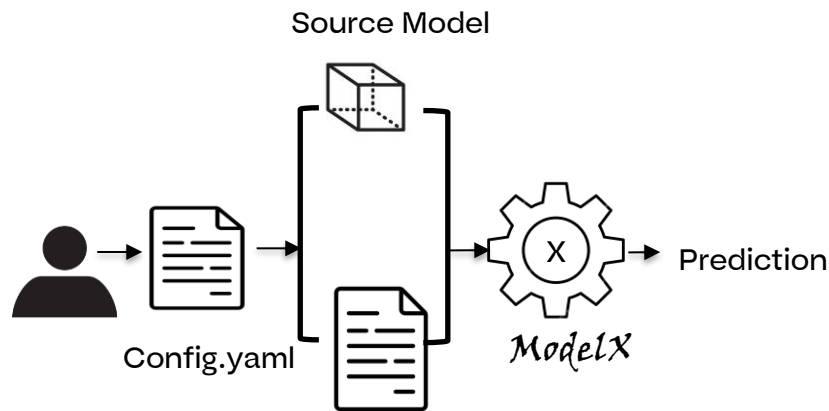
- Design a distance measure to quantify the “difficulty” of transferring knowledge between two datasets
- This measure can explain why SOTA does not work, and when different components of our solution is necessary or sufficient



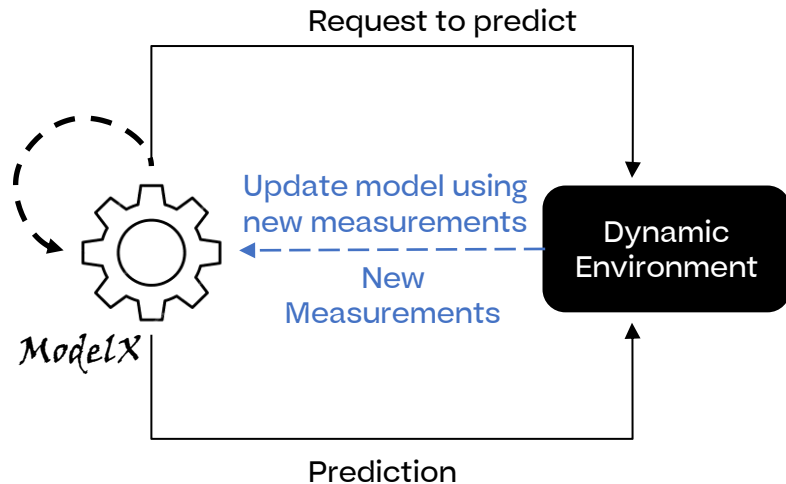
UNFINISHED

Distance between subspaces using
Grassmanian manifold

ModelX Can be Use in both Online and Offline Scenarios



Scenario 1: *ModelX* can be used during offline scenarios



Scenario 2: *ModelX* can be used during online scenarios

Experimental Setup



Metrics

Mean Squared Error



Datasets

11 HPC and 4 machine learning datasets



Scenarios

Cross applications
Cross architectures
Online job scheduling

App.	Example Features
HPC Datasets	
CG, LU, FT, Kripke1, CoMD, N=3180, p=7	Power Cap, Task Count, Core Count, Placement, and Bandwidth, runtime
Kripke2, N=17386, p=23	DRAMPowerPerNode, ProcessorPowerPerNode, Ranks, App-specific input parameters, OMP, PKG_LIMIT, DRAM_LIMIT, AvgInst, AvgIpc, AvgArithFpu, AvgFreq, AvgTemp, ProcessorPower, DRAMPower, Nesting Order, ExecTime
Hypre, N=50396, p=21	DRAMPowerPerNode, ProcessorPowerPerNode, Ranks, OMP, PMX, NS, MU, AvgIPC, Smoother, AvgTSC, AvgTemp, ProcessorPower, DRAMPower, Solver-related parameters, ExecTime
XSbench and OpenMC on SB, N=200, p=121	NumThread, InputSize, EfficiencyLoss , perf::[MEM DTLB LLC]_[MISS STALL], perf::[L1 L2 L3]_[MISS STALL]
XSbench and OpenMC on BGQ, N=200, p=145	NumThread, InputSize, Efficiency-Loss , PEVT_[XU AXU L1P STL]_[MISS], PAPI_[BR STL SYC]_[STALL CYC MISS]
ML Datasets	
Airfoil, N=1503, p=5	Frequency, Angle of attack, Chord length, Free-stream velocity, Suction side, Scaled sound pressure level
NO2, N=500, p=8	NO2 , Cars per hour, temperature, wind speed, temperature difference, wind direction, hour of day, day number
Crime, N=1994, p=127	population, householdsize, PctEmploy, PctIlleg, medIncome, perCapInc, PctPopUnderPov, ViolentCrimesPerPop
SkillCraft, N=3395, p=16	GameID, APM, SelectByHotkeys, AssignToHotkeys, MinimapRightClicks, NumberOfPACs, ActionLatency

ModelX Improves Prediction Accuracy Across Applications 93.5%

- Domain Divergence and Disjoint Features
- Overhead: Average test time adaptation 45.83s, average inference time 0.78s per query
- Number of features between 7 and 21
- For heterogeneous cases, *ModelX* has been compared against a supervised model with 100x data

Source	Target	Severity	Improvement	Winner	Best of Others
Airfoil	Airfoil	0	-82%	Source	Source
CoMD, CG,FT	Kripke[1]	0.5	56%	<i>ModelX</i> (Input Alignment)	Linear Probing
CoMD, FT. Kripke[1]	CG	0.65	69%	<i>ModelX</i> (Input Alignment)	Linear Probing
ComD, CG, Kripke[1]	FT	0.66	60%	<i>ModelX</i> (Input Alignment)	Fine Tuning
ComD, CG, LU, FT	Kripke[2]	0.60	95%	<i>ModelX</i> (Residual augmentation)	X
Kripke[2]	Hypre	0.70	99%	<i>ModelX</i> (Residual augmentation)	X
Hypre	Kripke[2]	0.70	99%	<i>ModelX</i> (Residual augmentation)	X

Homogeneous datasets

Heterogeneous

ModelX Improves Prediction Accuracy by 77% Across Architectures Compared to the Oracle using only 1-5% of the data

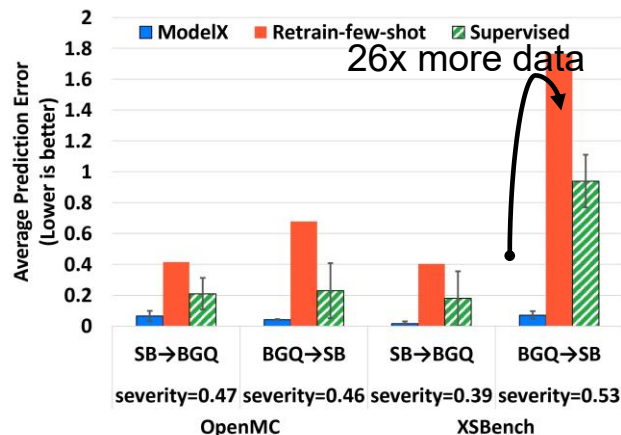
UNFINISHED

IBM BGQ – 143 features

B	C	E	F	G	H	J
PEVT_XU_BR_MISPRED_C	PEVT_LSU_ST_	PEVT_IU_AXU	PAPL_INT_INS	PEVT_LSU_LD_LA	PEVT_INST_QFP	PEVT_INST_XU
363275749.5	0	1.6723E+11	12205190201	11572	2449262430	2.5
181640446.8	6.625	8.3415E+10	6102393911	35687.5	1224932300	0.75
121043941.1	0	5.5493E+10	4065870710	18142.91667	816190208.8	0.166666667
90747775.13	488.9375	4.1512E+10	3048010431	3095.5625	611864963.6	0.125
73443709.1	542.3	3.3557E+10	2438308036	4484274.3	489484843.8	0.1
61696261.88	860.5416667	2.8135E+10	2032515138	6241053.542	408034358.4	0.083333333
53169009.32	1114.214286	2.4245E+10	1741726499	6890697.179	349656812.4	0.071428571

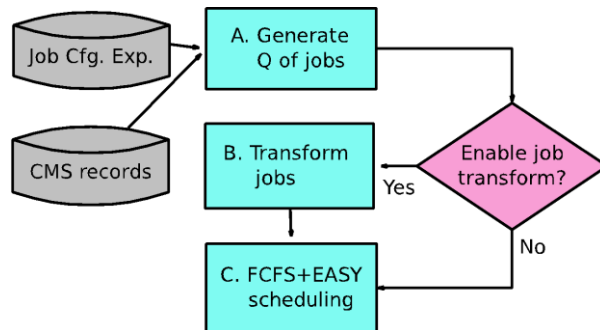
Intel Sandy Bridge – 121 features

B	C	D	E	F	G
MEM_LOAD_RETI	PAPL_L2_TCA	perf::INSTRUCTIONS	perf::NODE-STORI	perf::DTLB-STORE-M	INSTS_WRITTEN
2576406999	1577467923	24327658760	2	38137	103802174
1276757083	801386841	13069615784	5246383	10682.5	61142597
852924026.3	535450194.3	8700637454	2671912.333	4611.666667	56508432.67
635040848.8	401738887.3	6521345416	2300212	2898	40303665.5

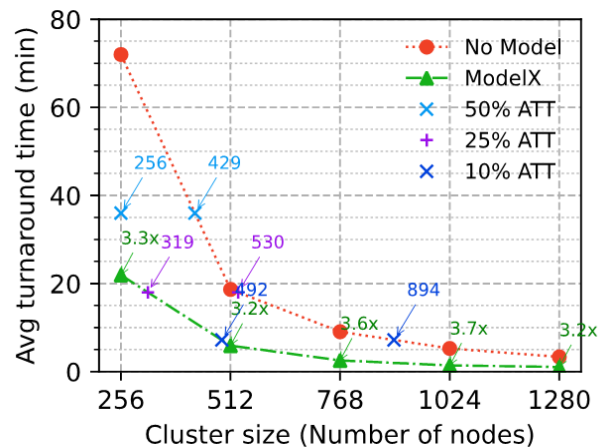


ModelX Reduces Average Turnaround Time by 71%

- Based on real-world job logs and performance measurement data from 6 HPC proxy applications
- Assumption: Jobs can run with a modified number of nodes than requested
- Scheduler asks *ModelX* to predict the execution time of a job using lesser number of nodes
- 3.4x time shorter turnaround time
- The state-of-the-practice scheduling method can perform as good by using up to 55% more nodes per job



- The Lassen job logs* collected over 2.5 year
- Extracted 70K jobs → 1-week's worth job
- Use the statistics of that week's job to create a stream of jobs



Summary