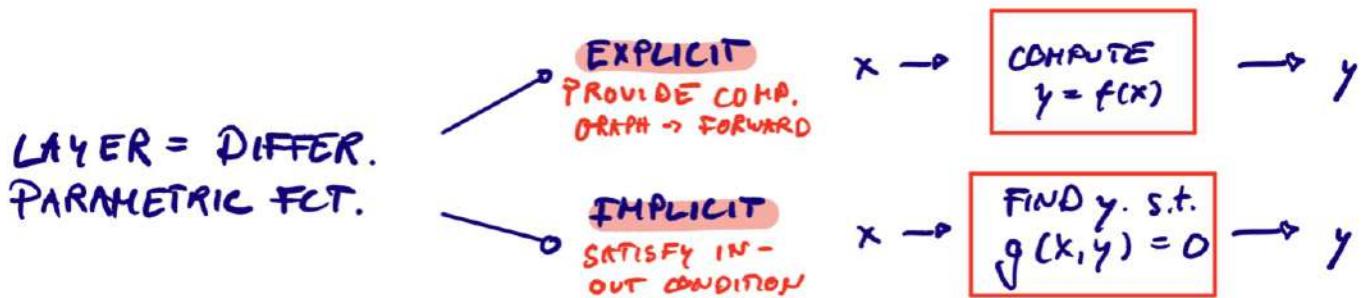


DEEP IMPLICIT LAYERS: NEURAL ODES, EQUILIBRIUM MODELS, ...

D. DUVENAUD (UoT), J. Z. KOLTER (CMU), M. JOHNSON (BRAN)



WHY USE IMPLICIT LAYERS?

1. POWERFUL REPRESENTATION
2. MEMORY EFFICIENCY
3. DESIGN SIMPLICITY
4. ABSTRACTION: "WHAT" - "HOW"

① DEEP EQ. M. ② NEURAL ODES ③ DIFF. OPTIM.

IMAGE CLASS.
SEMANTIC SEG.
LANGUAGE M.

CONT. TIME
SYSTEMS
GEN. MODELS
SMOOTH DENSITY
EST.

CONSTRAINED
OPTIMIZ.

HISTORICAL BACKGROUND

- RECURRENT BACKPROP
[PINEDA 87', ALMEIDA 87']
- APPLIED ENGINEERING
[RICO-MARTINEZ ET AL. 92', 95']

- DIFFERENTIABLE OPTIMIZATION

- \hookrightarrow STRUCTURED VAEs [JOHNSON ET AL. 16']
- \hookrightarrow D. DECLARATIVE NETS [GLOU ET AL. 16', 19']
- \hookrightarrow OPTNET [AMOS + KOLTER 17']
- \hookrightarrow CVXPY LAYERS [AGARWAL ET AL. 18']
- \hookrightarrow SAGNET [WANG ET AL. 18']
- \hookrightarrow SUBMODULAR OPTIM. [D'OLONGA + KRAUSE, 19']

SOLVING ODES \leftrightarrow IMPLICIT LAYER

$$\frac{dz}{dt} = f(z(t), t, \theta)$$

$$\Rightarrow z(t_1) = z(t_0) + \int_{t_0}^{t_1} f(z(t), t, \theta) dt$$

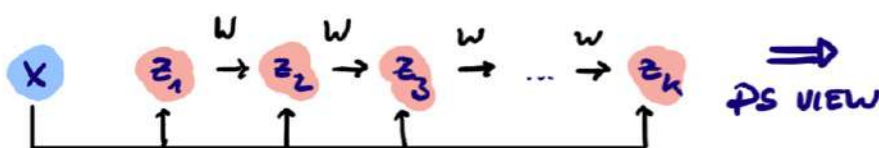
$$\Rightarrow y = \text{odeint}(f, x, t_0, t_1, dt)$$

\Rightarrow DROP-IN REPLACEMENT RESNET

- CONTINUOUS-TIME PHYSICS MODELS
 - \hookrightarrow INCORPORATE KNOWN STRUCTURES!
 - \hookrightarrow HAMILTONIANS / LAGRANGIANS

- CONTINUOUS NORMALIZING FLOWS
 - \hookrightarrow EASIER CHANGE OF VARS. COMPUTATION
 - \hookrightarrow ADAPTATION FOR DENSITY COMPUTATION
 - \hookrightarrow HOMEOMORPHISMS ON POINT CLOUDS

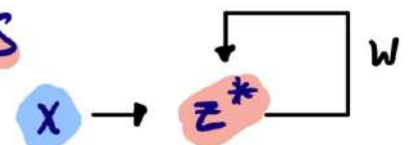
⑥ MATHEMATICS OF IMPLICIT LAYERS



WEIGHT-TYING
+ RE-INJECT!

$$z_{i+1} = \sigma(W z_i + x)$$

1 of 3



- \hookrightarrow NET CONVERGES TO FIXED / EQUILIBRIUM POINT
- \hookrightarrow HOW TO BACKPROP THROUGH? \Rightarrow MEMORY PROBLEM!

SOLUTION: IMPLICIT FUNCTION THEOREM

$f: \mathbb{R}^p \rightarrow \mathbb{R}^n, a_0 \in \mathbb{R}^p, z_0 \in \mathbb{R}^n$ s.t. w.r.t. z

1. $f(a_0, z_0) = 0$

2. f CONT. DIFF. W. NON-SING.

JACOBIAN $\partial_z f(a_0, z_0)$

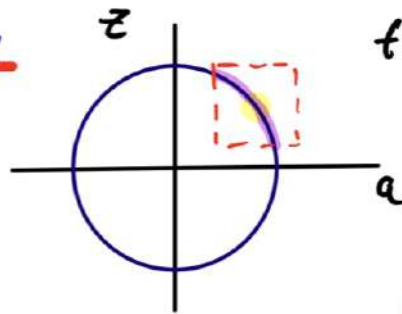
$\Rightarrow \exists$ OPEN SET $S_{a_0} \subset \mathbb{R}^p, S_{z_0} \subset \mathbb{R}^n$
CONTAINING a_0, z_0 & UNIQUE CONT.

FCT: $z^*: S_{a_0} \rightarrow S_{z_0}$ S.T.

1. $z_0 = z^*(a_0),$

2. $f(a, z^*(a)) = 0 \forall a \in S_{a_0}$

3. z^* DIFF. ON S_{a_0}



$$f(a, z) = a^2 + z^2 - 1 = 0$$

LOCAL SOLUTION MAPPING FCT.

\hookrightarrow FROM EXISTENCE TO DERIVATIVE EXPRESSION:

$$f(a, z^*(a)) = 0, \forall a \in S_{a_0}$$

$$\partial_a f(a, z^*(a)) + \partial_z f(a, z^*(a)) \partial_z z^*(a) = 0, \forall a \in S_{a_0}$$

$$\Rightarrow \partial_z z^*(a_0) = -[\partial_z f(a_0, z_0)]^{-1} \partial_a f(a_0, z_0)$$

CAN DIFFERENTIATE SOLUTION MAPPING FCT. WRT. PARAMS BY EVALUATING-
DERIV. OF f CLOSE TO
SOLUTION POINT!

$$\Rightarrow \text{FIXED POINT VERSION: } \partial_z z^*(a_0) = [I - \partial_a f(z_0, a_0)]^{-1} \partial_a f(z_0, a_0)$$

AUTODIFF. CONNECTION

1. JCP/FORWARD: $v \mapsto \partial f(x) v$

2. VJP/REVERSE: $w \mapsto w^T \partial f(x)$

$$w^T \partial_z z^*(a_0) = u^T f(z_0, a_0)$$

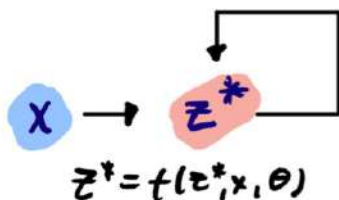
$$\text{with } u^T = w^T + u^T \partial_a f(z_0, a_0)$$

\Rightarrow BACKWARD PASS \Leftrightarrow LINEAR FIXED POINT IN TERMS OF VJPs!

① DEEP EQUILIBRIUM MODELS [BAI ET AL., 19] - DEQ

\Rightarrow REPLACE ENTIRE NET BY EQ. LAYER: $z^* = f(z^*, x, \theta)$

\hookrightarrow USE STABLE FP SOLVER ALGORITHM



BACKWARD VIA
IMPLICIT GRADS:

$$\partial_l \ell(\theta) = \partial_\theta \ell(z^*, y)$$

$$\cdot [I - \partial_a f(z^*, x, \theta)]^{-1}$$

$$\cdot \partial_z f(z^*, x, \theta)$$

\rightarrow PRACTICAL DETAILS: ANDERSON ACCELERATION

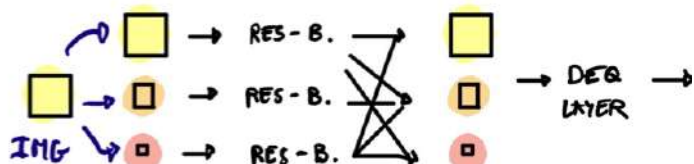
FOR FP ITERATIONS \rightarrow FORWARD + BACKWARD!

\rightarrow THEORY RESULT I: SINGLE-2. DEQ CAN REPRESENT ANY FFW DEEP NET.

\rightarrow THEORY RESULT II: SINGLE-L. DEQ CAN REPRESENT ANY MULTI-LAYER DEQ.

\rightarrow BUT: EXISTENCE EQ. POINT, UNIQUENESS, STABILITY?!

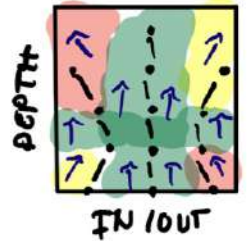
\rightarrow MULTI-SCALE
REPRESENTATIONS
 \hookrightarrow MIXING



1-2 of 3

② NEURAL ORDINARY DIFFERENTIAL EQ. [CHEN ET AL., 18]

$y = \text{odeint}(f, x, t_0, t_n, dt)$ \Rightarrow USE ANY SOLVER: EULER, RK, ADAPTIVE
 $\hookrightarrow f$ CONT. DIFF. + LIPSCHITZ \hookrightarrow FIT LOCAL POLYNOMIAL!



\hookrightarrow "PROBLEM": "INPUT PATHS" CAN'T OVERLAP! \rightarrow ONLY BIJECTIVE TRANSFORMS
 \hookrightarrow OFTEN DYNAMICS BECOME MORE COMPLEX W. TRAINING \rightarrow MORE FET. EVALS

$$\frac{\partial}{\partial t} \frac{\partial \mathcal{L}}{\partial z(t)} = \frac{\partial \mathcal{L}}{\partial z(t)} \frac{\partial f(z(t), \theta)}{\partial z} \Rightarrow \frac{\partial \mathcal{L}}{\partial \theta} = \int_{t_0}^{t_f} \frac{\partial \mathcal{L}}{\partial z(t)} \frac{\partial f(z(t), \theta)}{\partial \theta} dt$$

\hookrightarrow CONTINUOUS-TIME BP VIA JOINT METHOD: AUGMENT TRACE W. USP \rightarrow OCL1

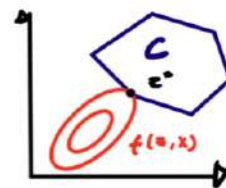
\hookrightarrow USE WHEN YOU CARE ABOUT TRAJECTORY \rightarrow CHANGE OF VAR. \rightarrow NORM. FLOWS

\hookrightarrow TRICK TO GET AROUND JACOBIAN TRACE: $\text{tr}(\pi) = E_{v \sim N(0,1)} [v^T \pi v]$
MUTCHINSON'S TRACE ESTIMATOR

\hookrightarrow NUMERICAL ERROR \Rightarrow PROPORTIONAL TO SOLVER TOLERANCE!

③ DIFFERENTIABLE OPTIMIZATION

LAYER: $z^* = \arg \min_{z \in C} f(z, x)$



\hookrightarrow SOLUTION TO CONSTRAINED OPT. \Leftrightarrow SOLUTION TO KKT CONDITIONS

\hookrightarrow VIEW OPTIM. PROCEDURE AS ITERATION \Rightarrow E.G. PROJECTED GD

$$z_{k+1} = \text{Proj}_C [z_k - \alpha \nabla f(z_k, x)] \rightarrow \text{AGAIN: IMPLICIT FET. THEOREM}$$

\hookrightarrow APPLICATIONS: LEARN CONVEX POLYTOPES, MNIST SDDUKO, HVAC MPC

\hookrightarrow CONVEX LAYERS \rightarrow DIFFERENTIABLE CONVEX MODELING \Rightarrow EMBED OPTIM. PROBLEM!

FUTURE DIRECTIONS AND OPEN PROBLEMS

DEQ

- DROP-IN INPUT REPLACEMENT
- SUPERVISED LEARNING
- STANDARD UNSUPERVISED

U. - ODE

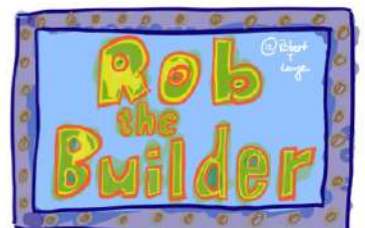
- CONT. TS \rightarrow IRREGULAR \rightarrow PHYSICS
- FLEXIBLE DENSITY
- HOMEOM.

??? REGULARIZING DEQS / NEURAL ODES FOR FASTER SOLVING \rightarrow PENALIZE DYNAMICS

??? ADAPT ARCHITECTURES TO EXPLOIT MEMORY ADVANTAGES \rightarrow DEQ "CELL" NAB

??? SCALING / APPLYING LATENT SDES

??? PDE SOLUTIONS AS A LAYER

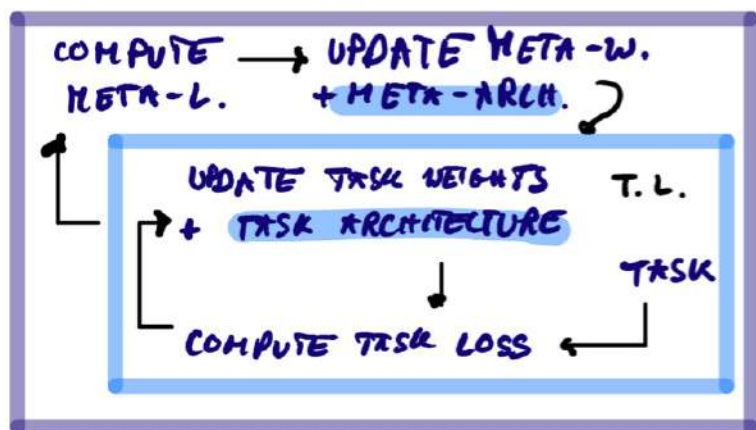


F. HUTTER (UO FREIBURG): 'META-LEARNING NEURAL ARCHITECTURES, INITIAL WEIGHTS, HYPERPARAMETERS AND ALGORITHM COMPONENTS'

① SAMPLE-EFFICIENT JOINT META-LEARNING OF MULTI-COMPONENTS

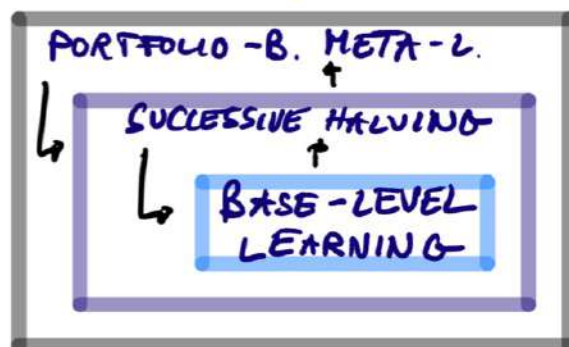
ELSSEN ET AL. 20' \Rightarrow META-NAS

\hookrightarrow SIMULTANEOUS WEIGHTS + ARCH.



ZIMMER ET AL. 20' \Rightarrow AUTO-P₁TORCH

\hookrightarrow MULTI-FIDELITY + ACROSS DATA



\hookrightarrow OFF-P. RL: SEARL, FRANKE ET AL. 20'

② META-LEARNING TO IMPROVE EXISTING ALGOS

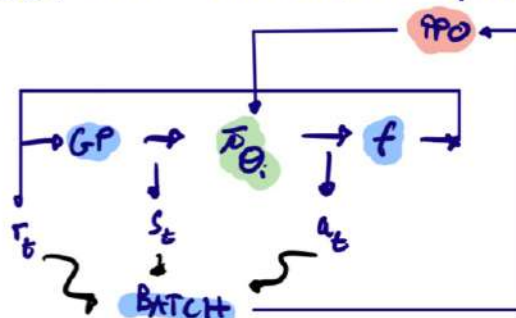
CHALLENGES WHEN LEARNING FROM SCRATCH:

GENERALIZATION TH. GUARANTEES SOTA PERF. \Rightarrow INSTEAD: IMPROVE EXISTING

\hookrightarrow ALSO: APPLICATION TO COMB. SOLVER ALGO SELECTION

\hookrightarrow ALSO: DYNAMIC ALGO CONFIG. \rightarrow ADAPT PARAMS TO CONTEXT

VDLPP ET AL. 20': META-L. ACQ. FCT. FOR BAYES OPT.



③ BENCHMARKS AS FOUNDATIONS OF MEASURABLE PROGRESS

HDP PLAYGROUND

RAZAVI ET AL. 20'

\hookrightarrow FAST TO RUN

\hookrightarrow MANY DIMS TO CONTROL

NAS - BENCH - XYZ

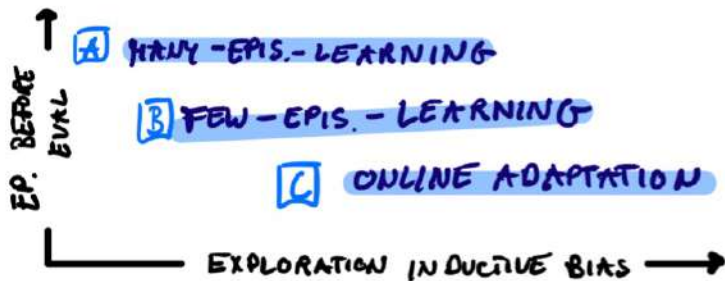
\hookrightarrow FAST EVALUATION BASED ON TABULAR LOOK-UPS \Rightarrow SINGLE RUN

\hookrightarrow SIEMIS ET AL. 2020 \Rightarrow SURROGATE MODELS

L. ZINTGRAF (UO OXFORD): 'EXPLORATION IN META-RL'

MOTIVATION: LEARNING A NEW TASK OFTEN REQUIRES EXPLORATION

① EXPLORE @ META-TEST TIME VS. EXPLORE @ META-TRAIN TIME



CHALLENGE: NEED TO EXPLORE META-TEST STRATEGIES AT META-TRAIN TIME

CAN BE HARD TO FIND!

STATE VALUE MAY VARY ACROSS TASKS

C. ONLINE ADAPTATION

→ EXPLORATION BONUS IN HYPER-STATE SPACE \Rightarrow HyperX [ZINTGRAF ET AL. 20'6]

$$\tau_{\text{env}} + \tau_{\text{hyper}} + \tau_{\text{error}}$$

EXTRINSIC REWARD

HYPER-STATE NOVELTY BONUS

BONUS FOR WRONG B. INF.

B. FEW-EPISODE-LEARNING

- MAX. RETURN AFTER N EPISODES
- GRADIENT-BASED ADAPTATION
- AGGREGATION-BASED ADAPT.

C. ONLINE ADAPTATION

- MAX. EXP. RETURN ONLINE
- SOLUTION VIA BAYES-ADAPTIVE AGENTS \Rightarrow DUFF & BARTO 02'
- TASK BELIEF \rightarrow OPTIMAL ACTION UNDER UNCERT.
- APPROX. INFERENCE - variBAD [ZINTGRAF ET AL 20'a]

BELIEF INFERENCE VIA VAE

BAYES-ADAPTIVE POLICY

$$b_t \leftarrow q(\mathbf{m} | \tau_{t:b})$$

$$\pi(a_t | s_t, b_t)$$

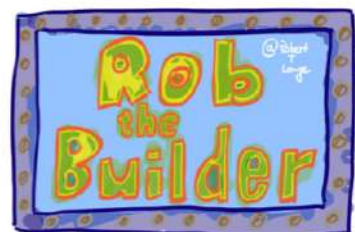
\Rightarrow PROBLEM w. HARD EXPLORATION!

B. FEW-EPISODE-LEARNING

- METACURE [ZHANG ET AL. 20']
- MAX. INFO GAIN OF EXPLORATION POLICY
- EXPLORE-THEN-EXECUTE [LIU ET AL. 20']
- EXPLICIT LEARNING OF EXPLORATION π TO RECOVER TASK ID EMBEDDINGS

* MANY-EPISODE-LEARNING

- OPEN Q! FUTURE RESEARCH
- HOW TO DEFINE + SEARCH?
- WOULD LIKE TO LEARN HOW TO SET EXPLORATION BONUSES!

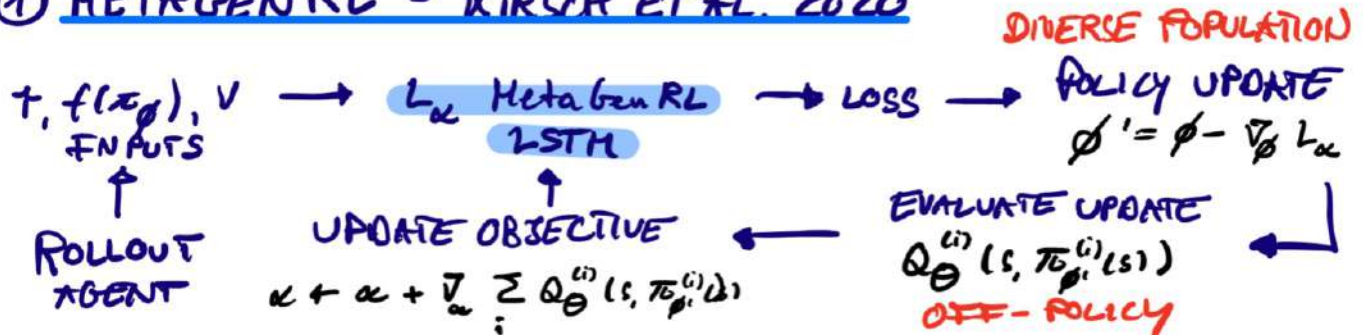


L. KIRSCH (IDSIA): 'GENERAL META-LEARNING'

MOTIVATION: MINIMIZE HAND-CRAFTED INDUCTIVE BIASES

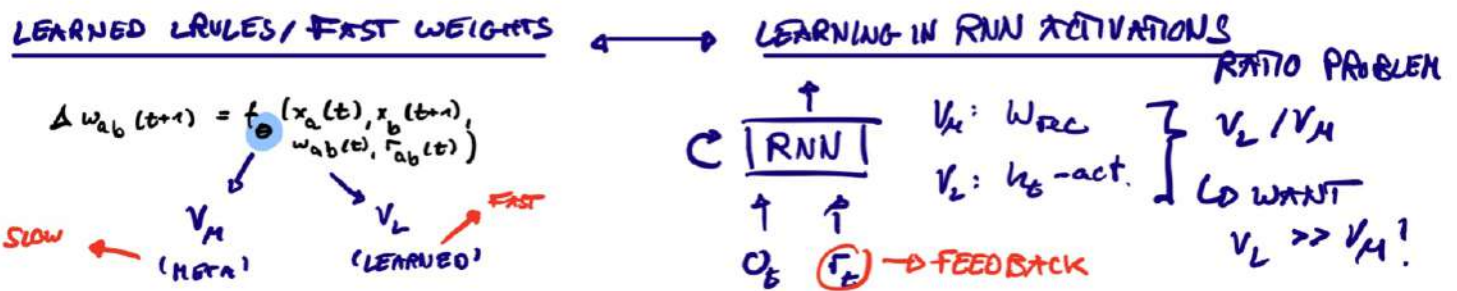
↳ NEED FOR BROAD GENERALIZATION TO BE APPLICABLE

① METAGENRL - KIRSCH ET AL. 2020

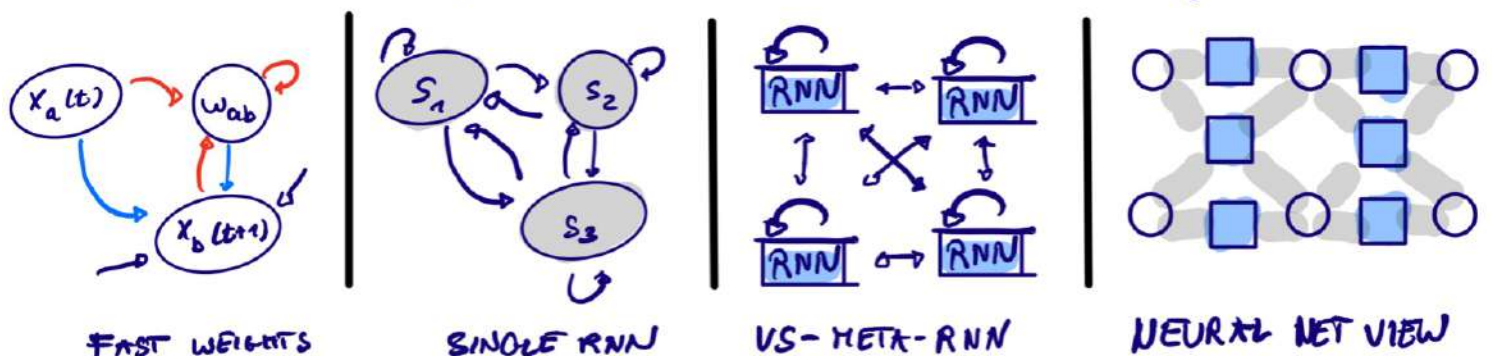


⇒ GENERALIZATION TO UNSEEN ENVS + OUTPERFORMING PPO VIA
 ⇒ EXTENSION BY LEARNED PG → OH ET AL 20' 2ND ORDER-GRADS!

② VARIABLE SHARED META-LEARNING (VS-M2)



↳ VS-M2: $s_{nnj} \leftarrow \sigma(b_j + \sum_i s_{nni} w_{ij} + \sum_i s_{nni} V_{ij})$ (labeled "INTERACTION")



↳ META-LEARN GENERAL LEARNING ALGO ⇒ MNIST vs. FASHION MNIST

③ BOOTSTRAPPING AI CONJECTURE

