## Lifelong learning paper

AIM: Build a bi-directional transfer learning system

DATA
↓
KNOWLEDGE
applicable to
past
AND future.



- ensembles rep$^n$ learnt independently on all tasks
- decides how to proceed on new data point
- 1 single decision making entity
  - ∴ combines all knowledge from all tasks

**Q** If the omni voter layer is deciding how to respond, how does it learn "future tasks"?

**Catastrophic forgetting**: Performance on old tasks falls rapidly as you learn new tasks

### Ways to overcome this

**Reallocate resources**
i.e, forget some things, learn new ones

**Add/build resources**
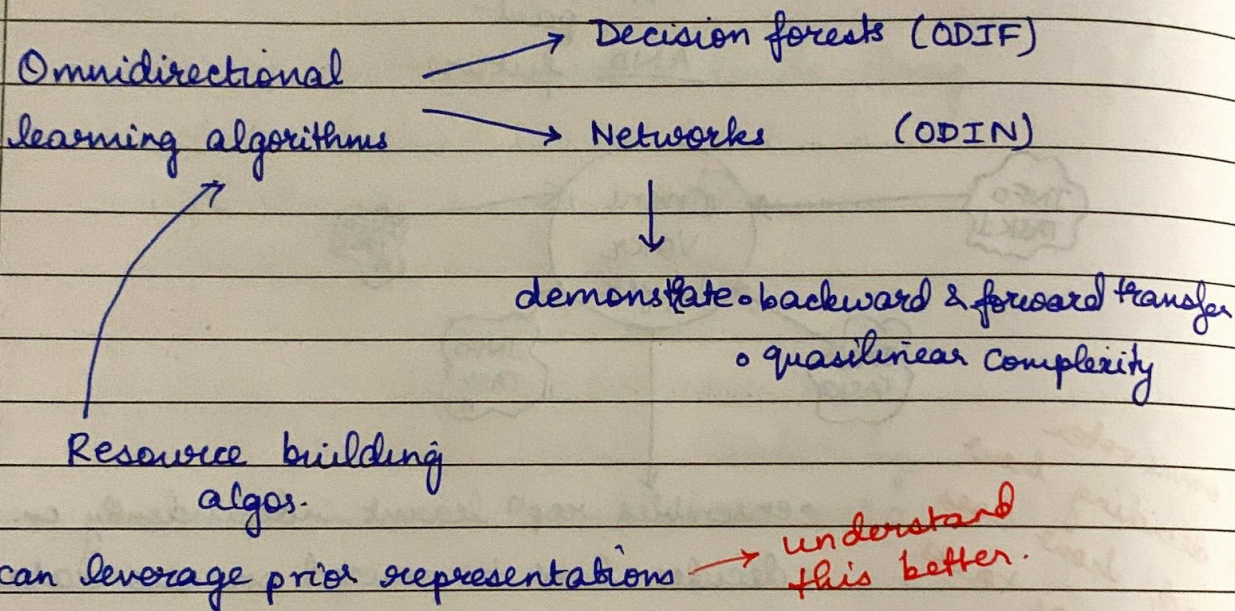i.e, add connec$^n$ ~~based~~ as part of development to make space for new info.

Neither exhibit omnidirectional learning
i.e, most cannot learn forward
& none learn backward

\* Read ProgNN paper → DeepMind.

\* What is quasilinear computational complexity?

Significant step forward:

Progressive NN

new tasks yield additional representational capac

Omnidirectional → Decision forests (ODIF)

learning algorithms → Networks (ODIN)

↓

demonstrate • backward & forward transfer

• quasilinear complexity

Resource building algos.

can leverage prior representations → understand this better.

Lifelong learning generalizes classical machine learning by:
    i) Environment of $T$ (possibly $\infty$) tasks instead of 1
    ii) Sequential data                instead of batch
    iii) Computational complexity constraints

         ∴ can't just retrain when new data

## GOAL IN LIFELONG LEARNING:

    Given new data & a new task,
         use all the existing data to achieve lower
         generalization error on all new & previous tasks

## IN CLASSICAL ONLINE SYSTEMS:

*Isn't this true even in lifelong learning?*

    Previously experienced tasks may recur,
         ∴ maintain / improve performance on those as well

NOTE: Only looking at task-aware scenarios
    ie, we know all task details for all tasks
    i.e, $h : X \times T \longrightarrow Y$

## COMPARED TO 9 REFERENCE ALGORITHMS

EVALUATION CRITERIA:

$$\text{Transfer efficiency} = \frac{\text{Gen}^n \text{ error of algo trained on task } t, \text{ data of } t}{\text{Gen}^n \text{ error of algo trained on all data}}$$

$R^t$ : risk of task $t$

$S_n^t$ : data assoc. with task $t$

$S_n$ : all data

$R^t(f(S_n^t))$ risk on task '$t$' by hypothesis learnt on ~~all data~~ task $t$ data

$R^t(f(S_n))$ risk on task '$t$' by hypothesis learnt on all data

$$\text{Transfer efficiency} = \frac{E(R^t(f(S_n^t)))}{E(R^t(f(S_n)))}$$

$$\uparrow$$
$$TE_n^t(f)$$

Algorithm '$f$' has transfer learned for task '$t$' with data '$S_n$' if & only if $TE_n^t(f) > 1$

<span style="float:left">task info known for all, but data appears sequentially</span> Respecting the 'streaming' nature of tasks:

Forward transfer efficiency = $\dfrac{\text{Risk of algo. with access to data from '}t'}{\text{Risk of algo. with access upto task '}t'.}$
$(FTE_n^t(f))$

& including last obs$^n$ of.

$$= \frac{E(R^t(f(S_n^t)))}{E(R^t(f(S_n^{<t})))}$$

An algo. has used data from prv. tasks to improve perf. on $t$ if FTE > 1

Backward transfer efficiency = $\dfrac{E(R^t(f(S_n^{<t})))}{E(R^t(f(S_n^t)))}$ = $\dfrac{\text{Risk with access upto '}t'}{\text{Risk with access to '}t'}$
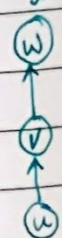$(BTE_n^t(f))$

An algo. has used data from future tasks to improve perf. on $t$ if BTE > 1

# Omnidirectional Algorithms

$$h(\cdot) = w \circ v \circ u(\cdot)$$

**for single learner:** (w → v → u)

hypotheses for
lifelong learning
can be decomposed
into 3 constituent parts

| decider | voter | representer |
|---|---|---|
| $\Delta y \to Y$ | $\tilde{x} \to \Delta y$ | $u : X \to \tilde{X}$ |
| produces a predicted label | maps transformed data to a posterior distribution | maps i/p $X$ to internal rep$^n$ space $\tilde{X}$ |
| | i.e, a score on response space $Y$. | |

## Step 1

'B' diff representers → one voter per representer
↓
B voters

Decider ensembles voters



## Step 2

Each voter ensembles representers → helps if each representer has learned complementary representations
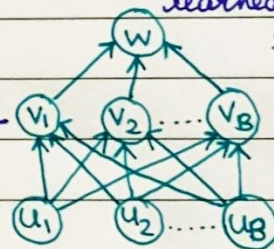
"Omni-voter layer" ←

ensembles all the existing representations, irrespective of the ~~order they were learnt in~~



Then, ensemble of voters feeds into the decider ~~helps if~~

* **Omnidirectional Forest**: Decision forest based instance of
  (ODIF)           ensembling representations

* **Omnidirectional Network**: Deep network based instance of
  (ODIN)           ensembling representations