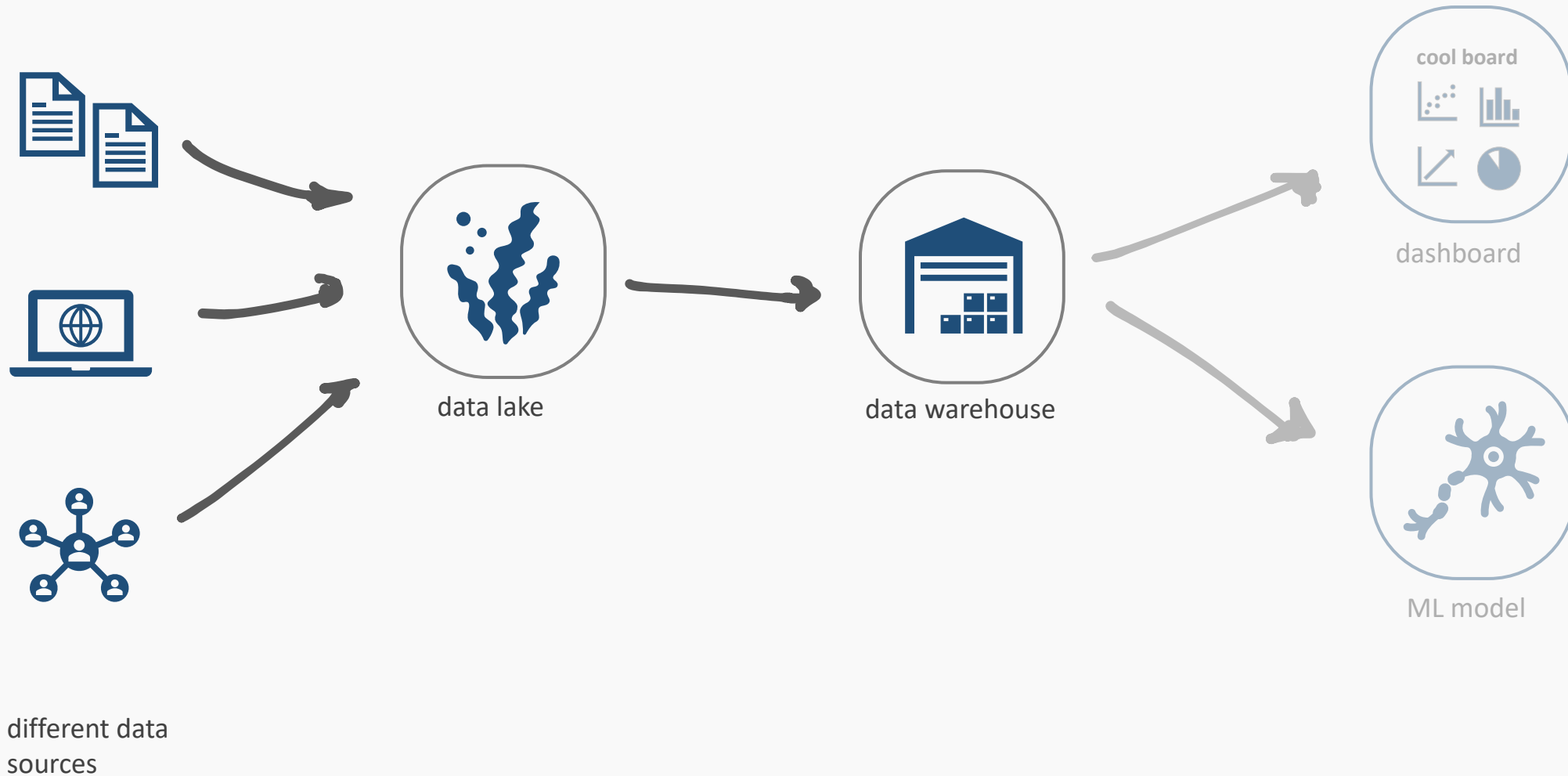kokchun giang
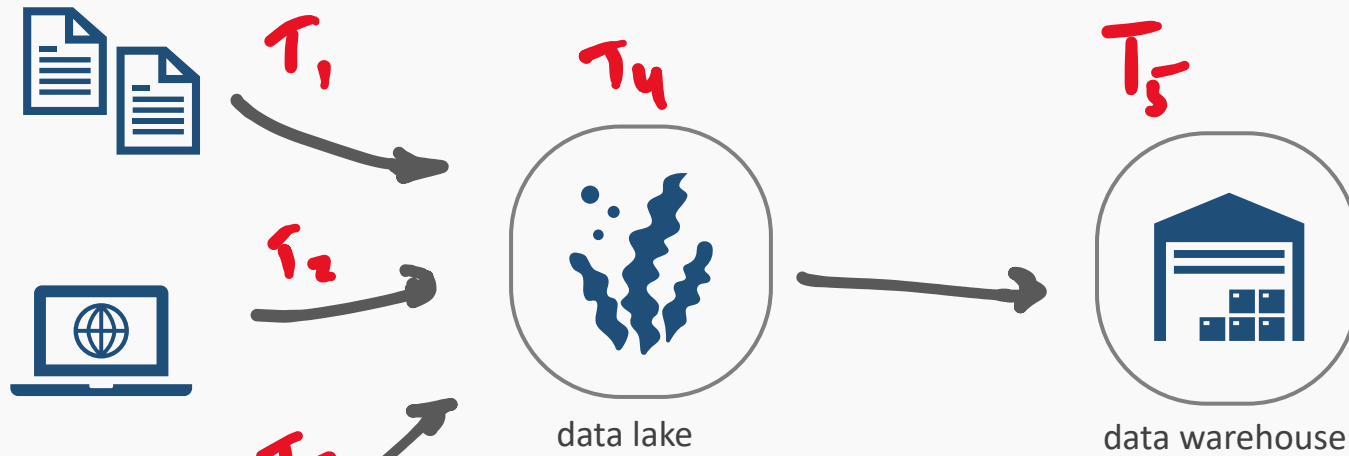
put your tasks in  directed acyclic graphs **(DAGs)** in **airflow** to create data pipelines

# let's take a look at a **data pipeline**



different data
sources

# let's take a look at a **data pipeline**



$T_1$ $T_2$ $T_3$ — extract data

data lake — $T_4$ — transform — unstructured raw data

data warehouse — $T_5$ — structured data ready for ML & dashboard

a how to solve this task ?

# use **cron** scheduling to solve it?
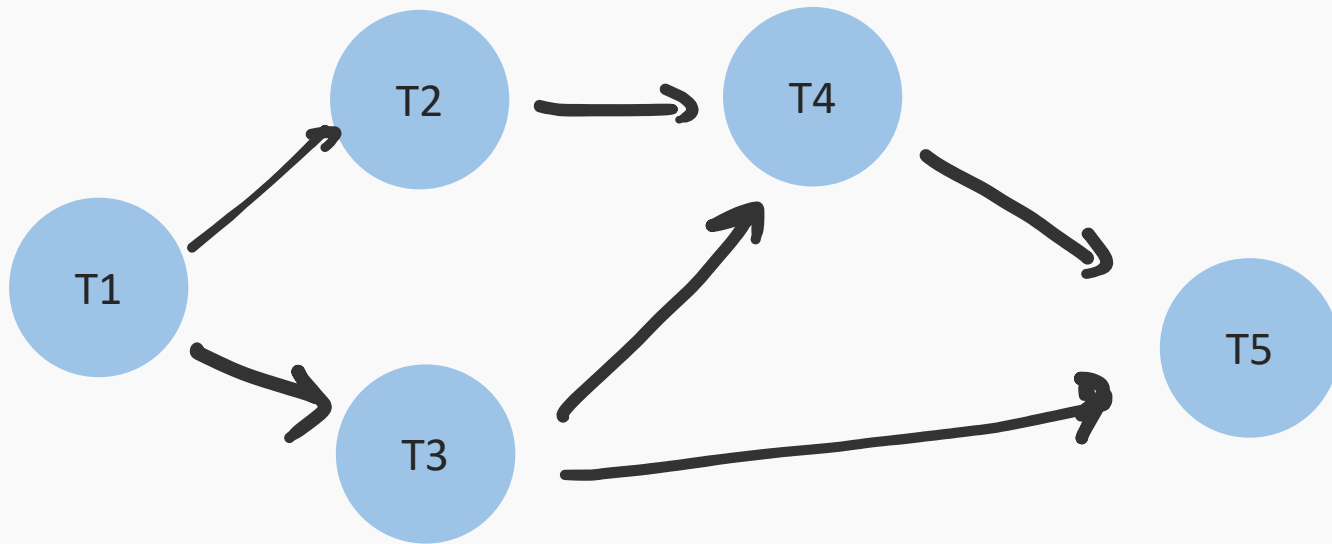


T1

T2 → ? T4 → T5

T3

a naive solution:

5 scripts $T_1 \rightarrow T_5$

× use cron to schedule
all scripts with 1h interval

× done ?

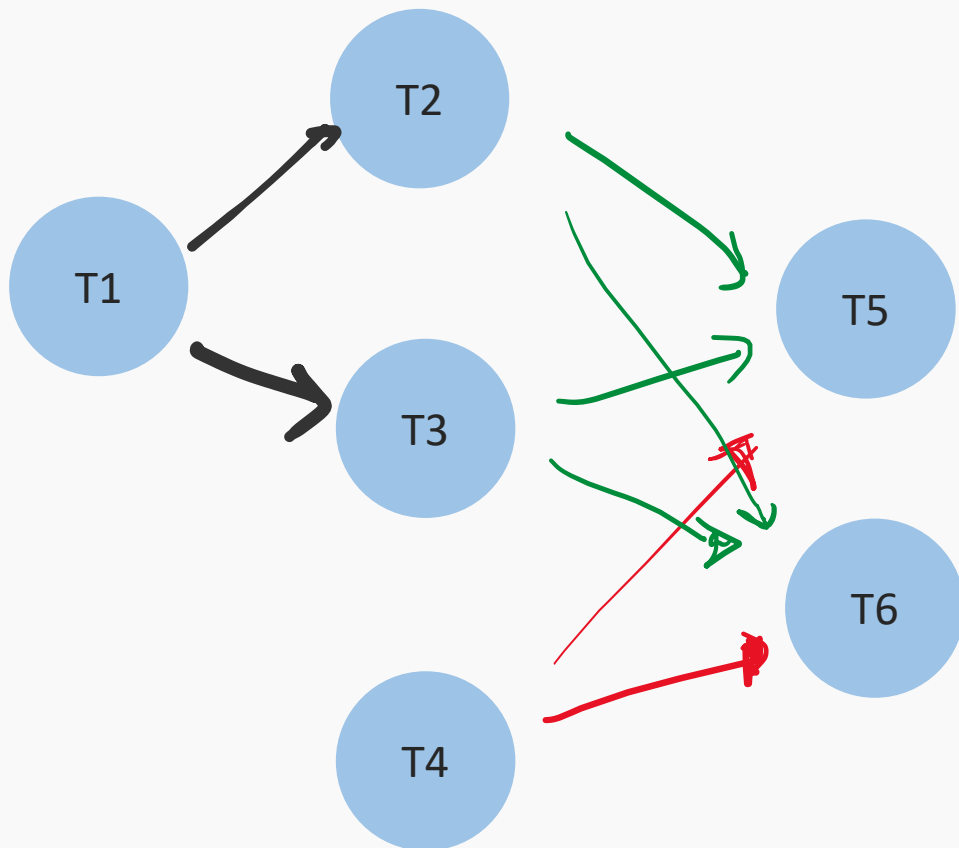# how about when we have **dependencies?**



a so $T_4$ waits for $T_2$ & $T_3$

* $T_5$ waits for $T_3$ & $T_4$

* can we schedule this with cron?
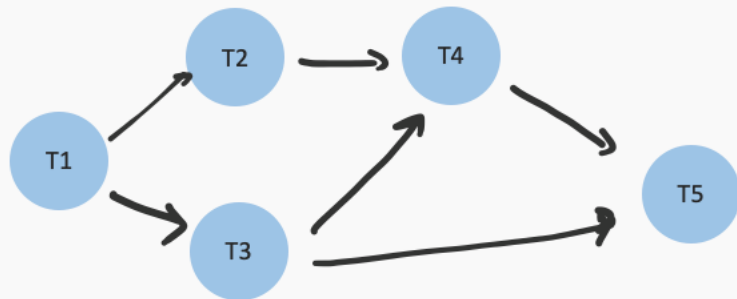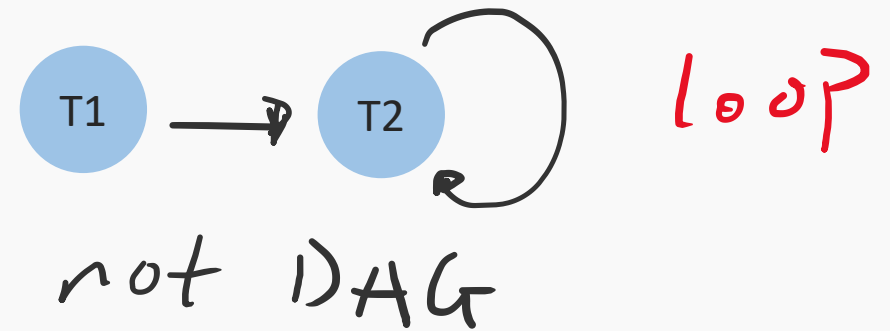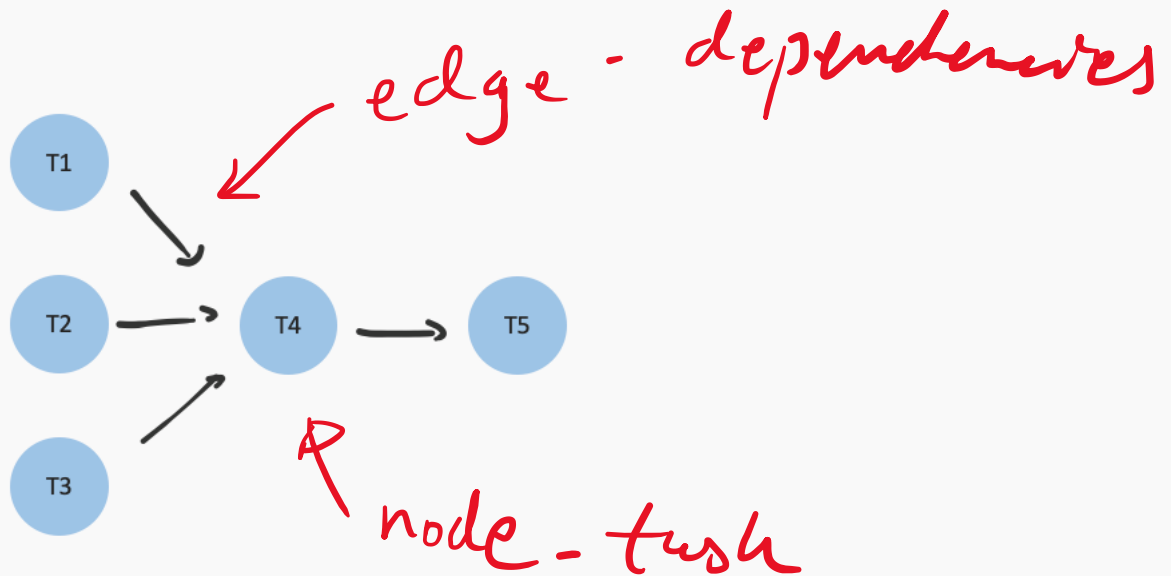
→ maybe but becomes complex!

# some more examples when **cron** is not enough



T5 should run when at least two of previous was successful

T5 should run when at least one of upstream succeeds, e.g. T6 is an email/discord/slack notification when some upstream tasks fails

# what are **DAGs**



edge - dependencies

node - task

T1 → T2 ⟲ loop

not DAG

T1 — T2 undirected

not DAG

# what are **DAGs**

**D**irected    $\forall$ tasks $\geq$ 1 upstream $\vee$ 1 downstr

**A**cyclic    tasks no dependency to itself, $\exists$ no loops

**G**raphs    relationship betw. tasks by nodes & vertices

DAGs are data pipelines with collection of tasks and relationsships between tasks
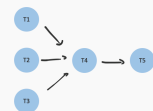
# orchestrate workflows with **airflow**

coordination and management of multiple tasks

define a task
(unit of work)

tasks organized
in DAGs

schedule the
execution of tasks

upstream/downstream
dependency

monitor progress in
Airflow UI

integrate with
other tools