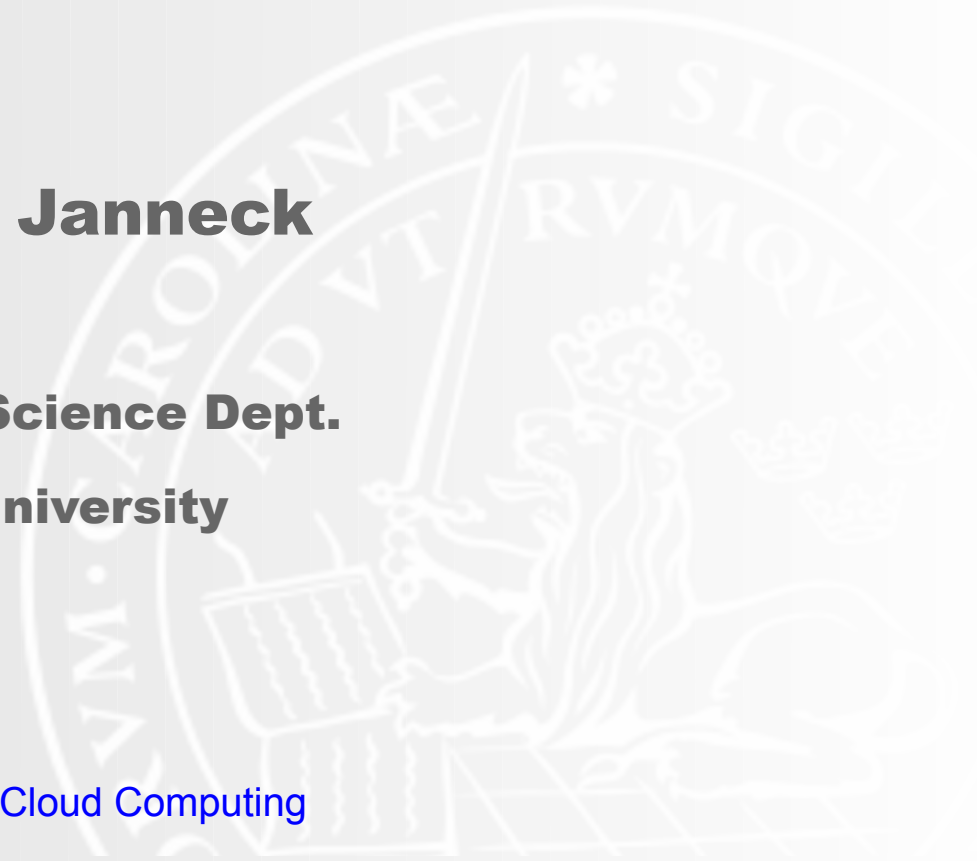# Distributed Computing II

## Jörn W. Janneck

### Computer Science Dept.

### Lund University

Introduction to Cloud Computing

# homework solution

One prisoner is chosen to be the counter.
He keeps a count, starting at 0.

If the counter goes into the room:
> If the switch is up:
>> Move it down.
>> Increment the count.
>> When the count reaches 99 (100–himself),
>> tell the warden everyone has been
>> in the room.
>
> If the switch is down:
>> Leave it.

If any other prisoner goes into the room:
> If the switch is down and he has not
> previously moved it up:
>> Move it up.
>
> Otherwise leave it.

**(Emma Fitzgerald)**

# consensus & data structures

Victor: Paxos

Manfred: DHT

# consensus

**P1** v1 ————————————— v3 →

**agreement**
   all correct processes end up with the same value
**termination**
   all correct processes will eventually make a decision

**P2** v2 ————————————— v3 →

**(strong) validity**
   the value decided upon is one of the input values
weak validity
   if all processes receive the same input value,
   all correct processes will decide on it

**P3** v3 ————————————— v3 →

**P4** v4 ————————————— X →

**P5** v5 ————————————— v3 →

# 2PC (two-phase commit)



phase 1
*propose & vote*

phase 2
*commit or abort*

v

**P1**

propose(v)

**P2**

**P3**

vote(Y|N)

commit()/abort()

**P4**

**P5**

agreement?
validity?
termination?

robustness?

# 3PC (three-phase commit)