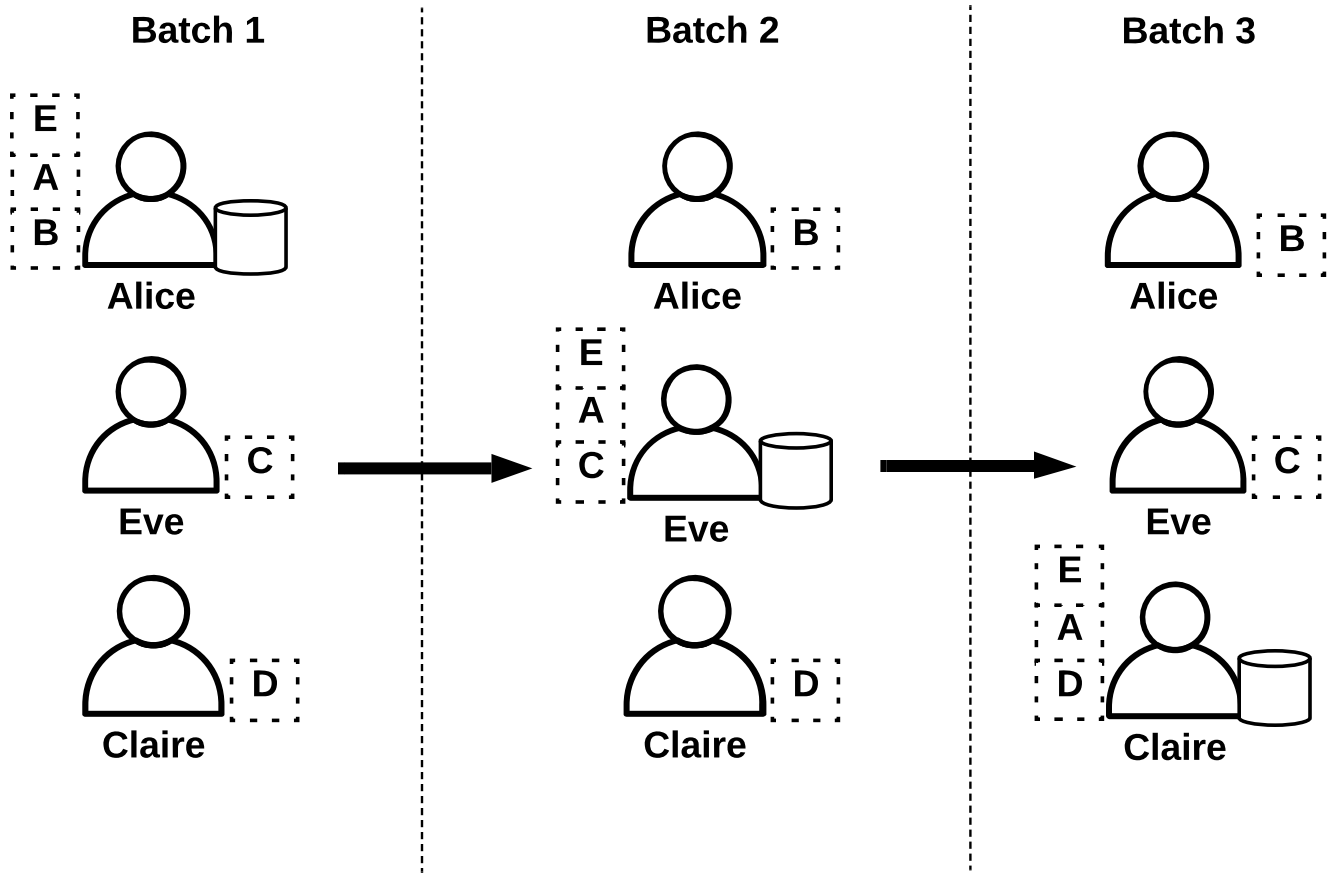
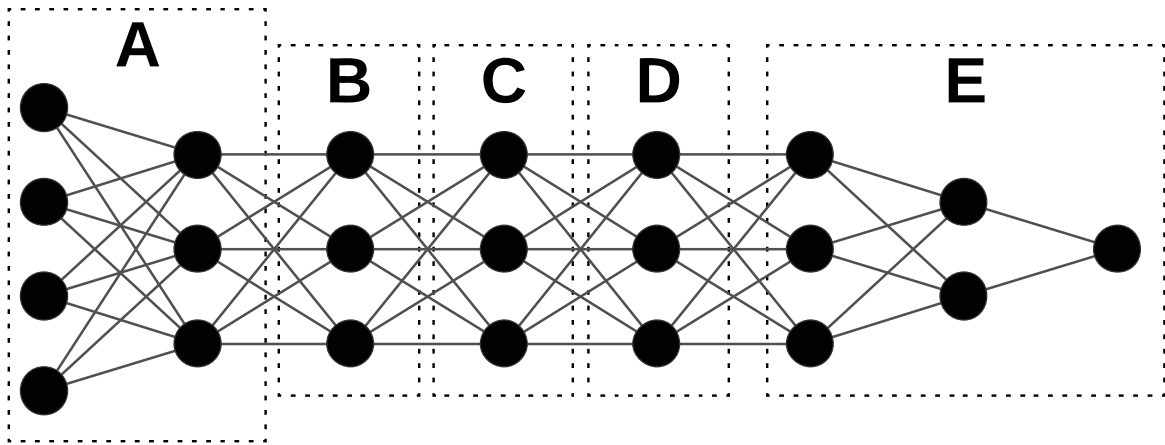
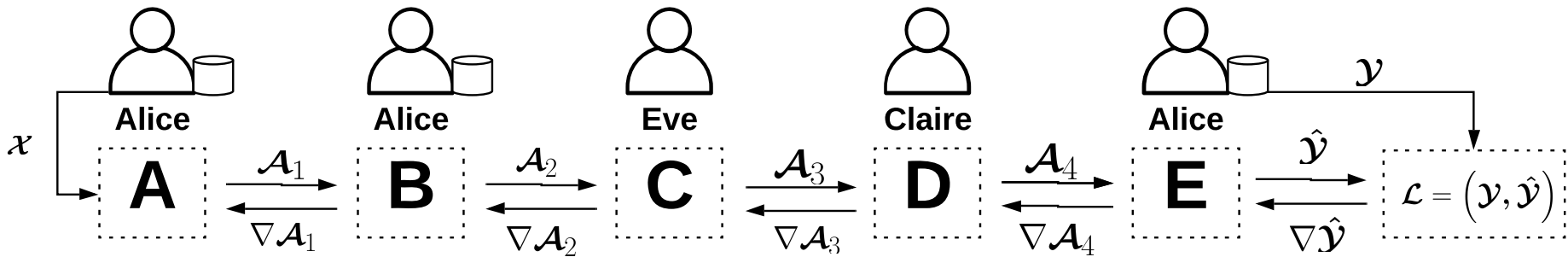


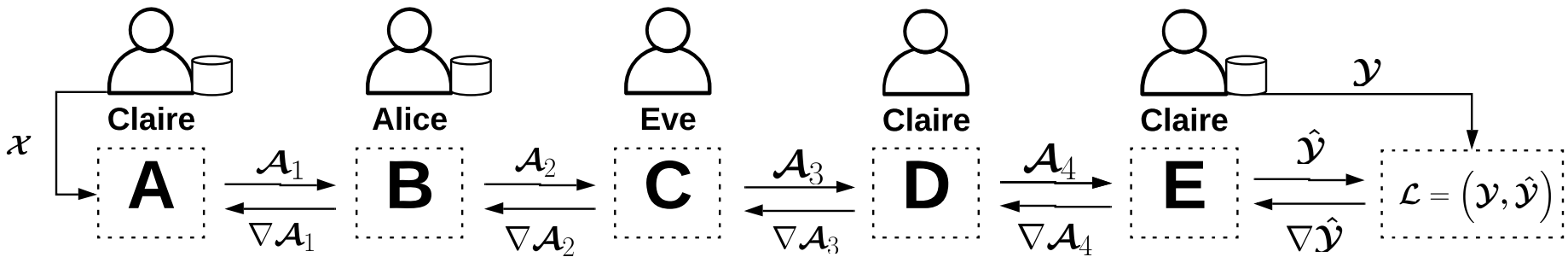
# Split Neural Network



1.



3.



# Vulnerability B

Eve can approximate the X values of Claire;

- Eve chooses  $X^\wedge$  values which approximately span the domain of  $X$
- Eve receives the activation signals produced by her  $X$  values
- Eve can statistically correlate the values of  $A$  that she received with the  $X^\wedge$  values that she entered. This correlation could perhaps be established using a second learning model with  $A$  as its inputs and  $X^\wedge$  as its label.
- When it's Claire's turn to train, Eve will receive  $A$  values from Claire.
- Using the model trained in the step before, Eve can make predictions as to the  $X$  values of Claire.

Key Takeaway;

Data owners should not also be involved in processing the model as this can leak information.

If Eve is the Researcher, Vulnerability A applies.

- Just replace  $M_0$  with  $B(A(X))$

+

2.

