

??  
 $t_{\text{axonomy.pdf}}$  Taxonomy of layer partitioning schemes across PEs.  
**Batch**  
 Partitioning  
 (Data-parallelism):  
?

Fmap  
 (Image)  
 Partitioning:  
 $11^3 \times 11^2$   
???

Output  
 Partitioning:  
 $K$   
 $K/p$

Input  
 Partitioning:  
 $C$   
 Importantly,  
 input partitioning requires an All-Reduce operation to aggregate partial results across PEs

??  
 Scheme Weight IFmap OFmap Reduction  
 All-Reduce Yes

??  
 hybrid partitioning scheme  
 fmap partitioning output partitioning

$$E_{\text{access}} = A_{\text{DRAM}} \times e \times (1 + \beta \cdot r)$$

$$(1) \quad \begin{matrix} e \\ \beta \\ r \\ A_{\text{DRAM}} \\ A_{\text{DRAM}} \\ A_{\text{DRAM}} \\ r \end{matrix}$$

INPP  
 All-Reduce Cost:

$$(2) \quad E_{\text{INPP}} = D_{\text{output}} \times \frac{2(p-1)}{p} \times e_{\text{comm}}$$