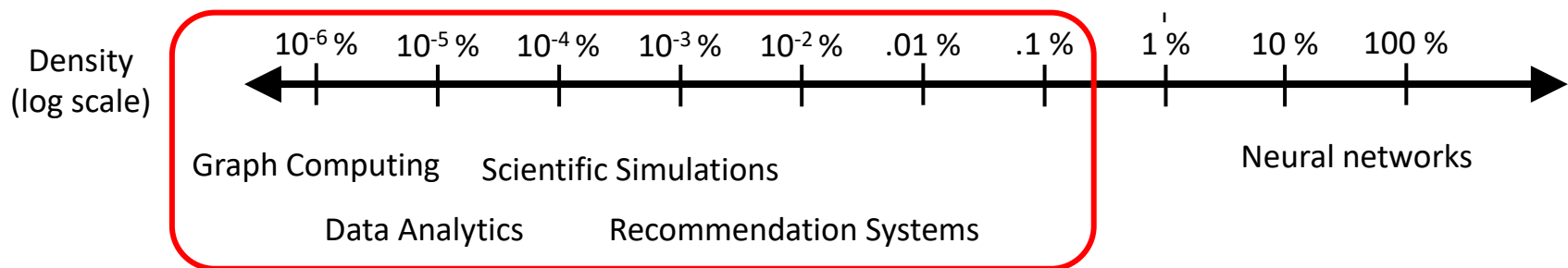


# Accelerating Sparse Tensor Algebra by Overbooking Buffer Occupancy

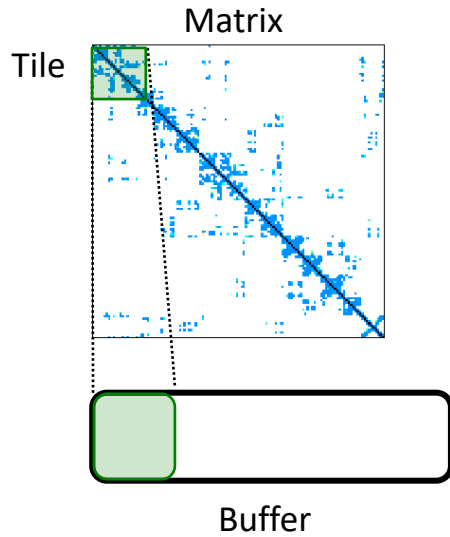
Fisher Xue, Nellie Wu, Joel Emer, Vivienne Sze



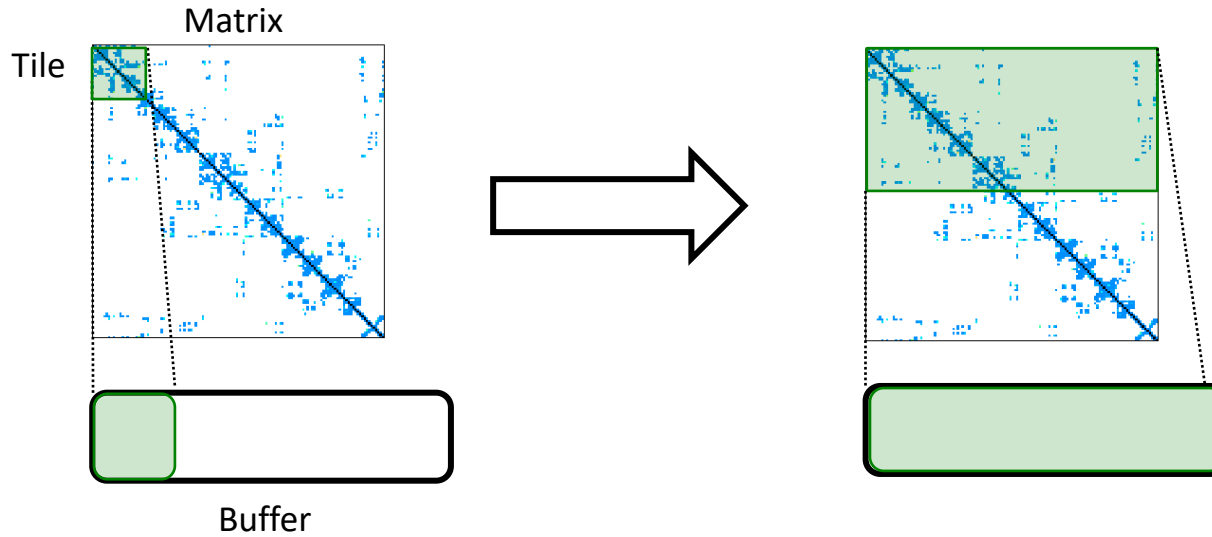
# Many applications operate on highly sparse tensors



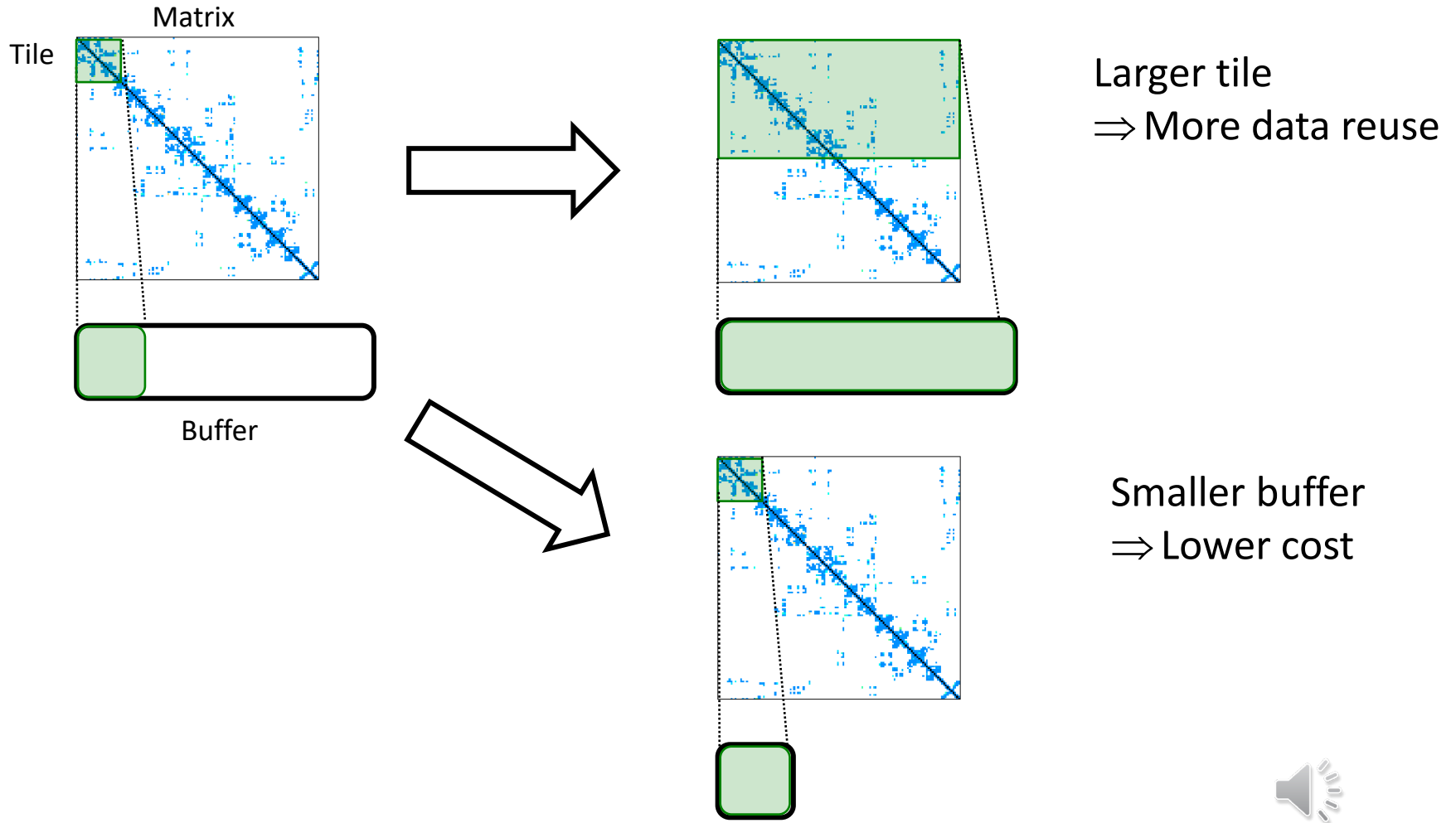
# Better tiling improves on-chip data reuse



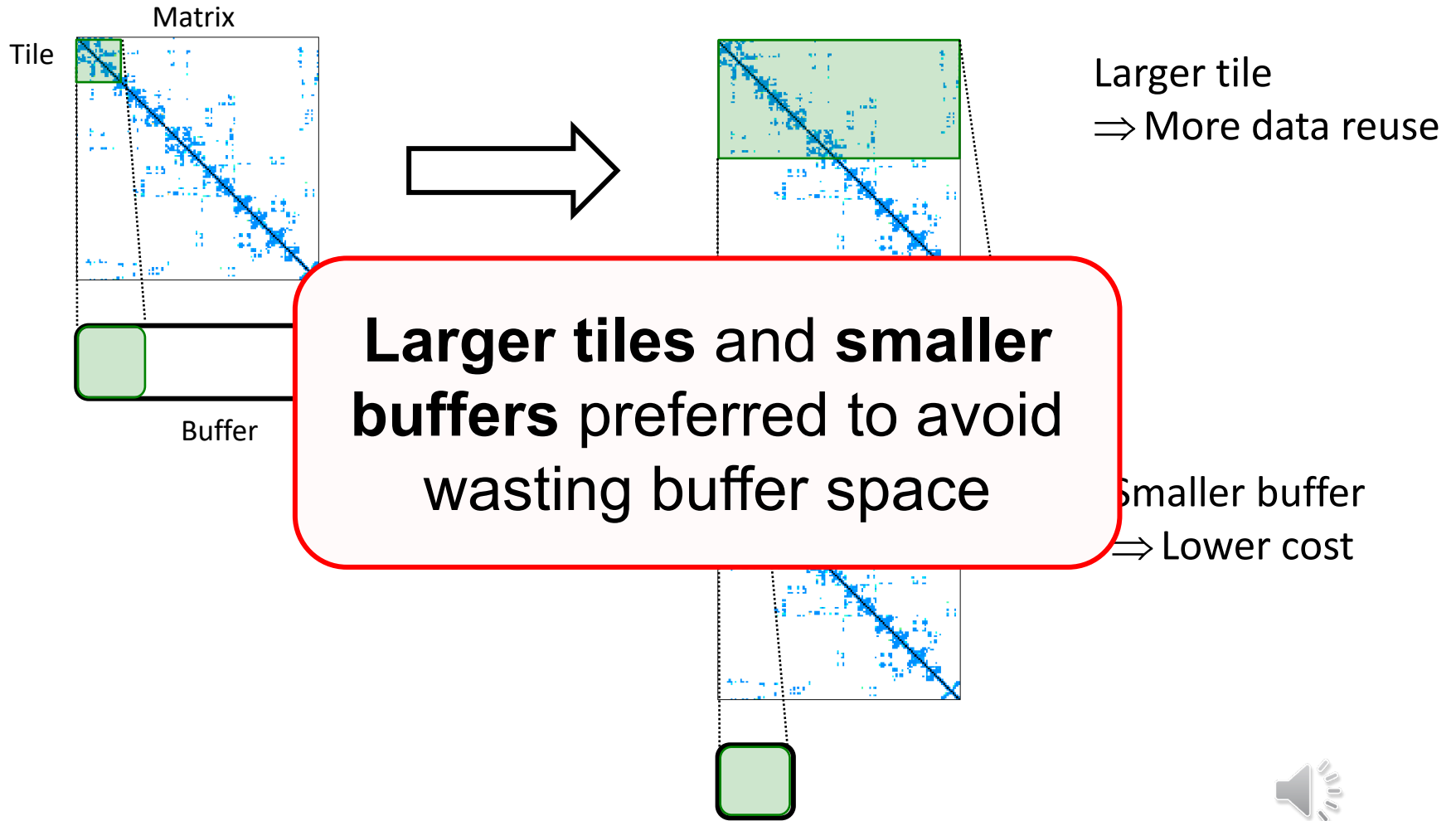
# Better tiling improves on-chip data reuse



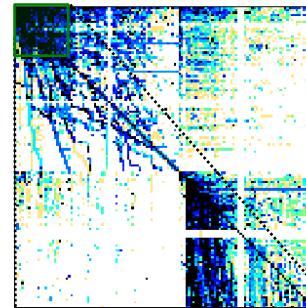
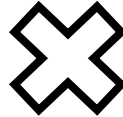
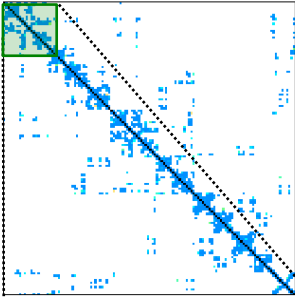
# Better tiling improves on-chip data reuse



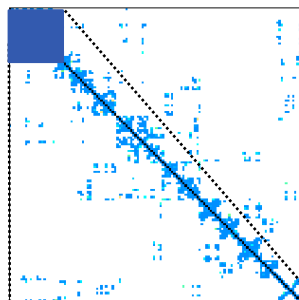
# Better tiling improves on-chip data reuse



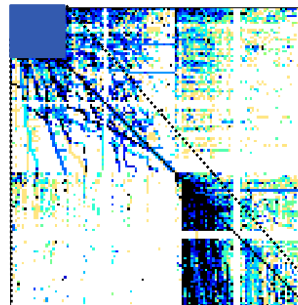
# Coordinate-space tiling



# Coordinate-space tiling



Nonzero

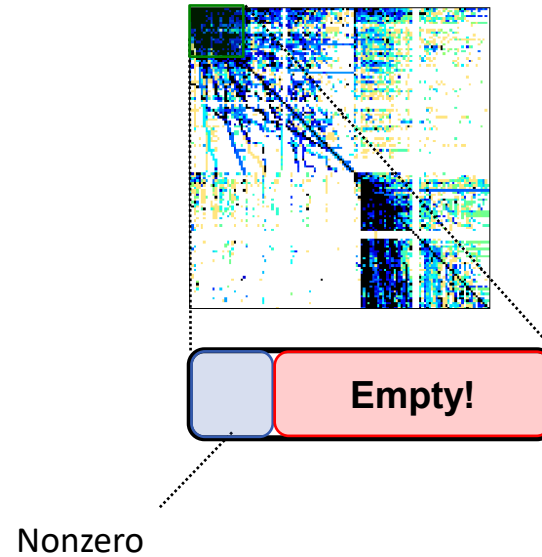
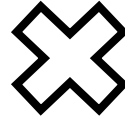
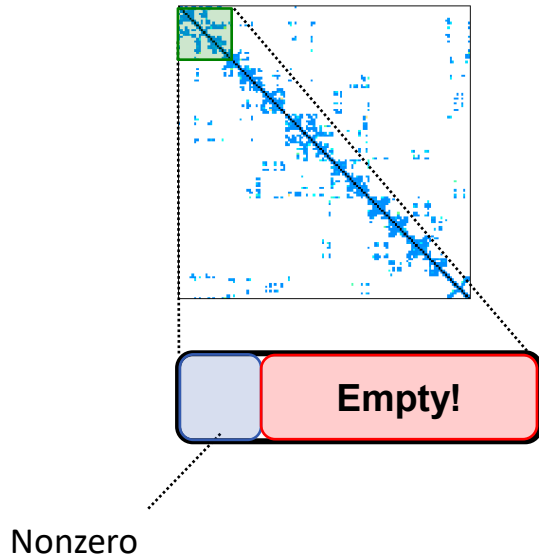


Nonzero

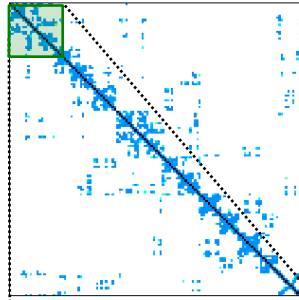




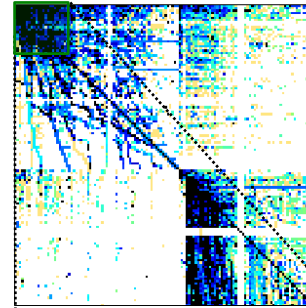
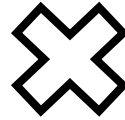
# Coordinate-space tiling



# Coordinate-space tiling

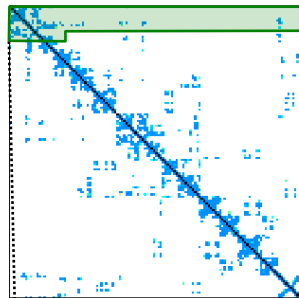


Empty!



Empty!

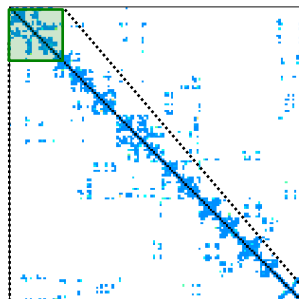
# Position-space tiling



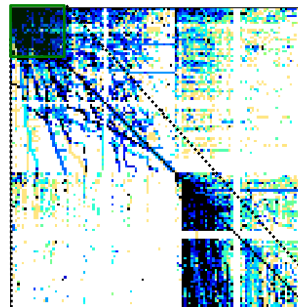
Nonzero



# Coordinate-space tiling

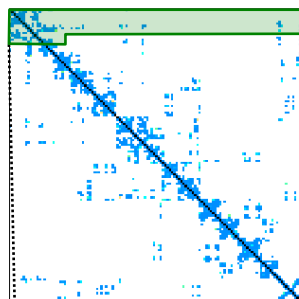


Empty!

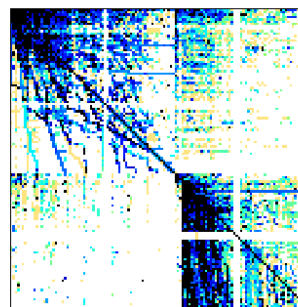


Empty!

# Position-space tiling



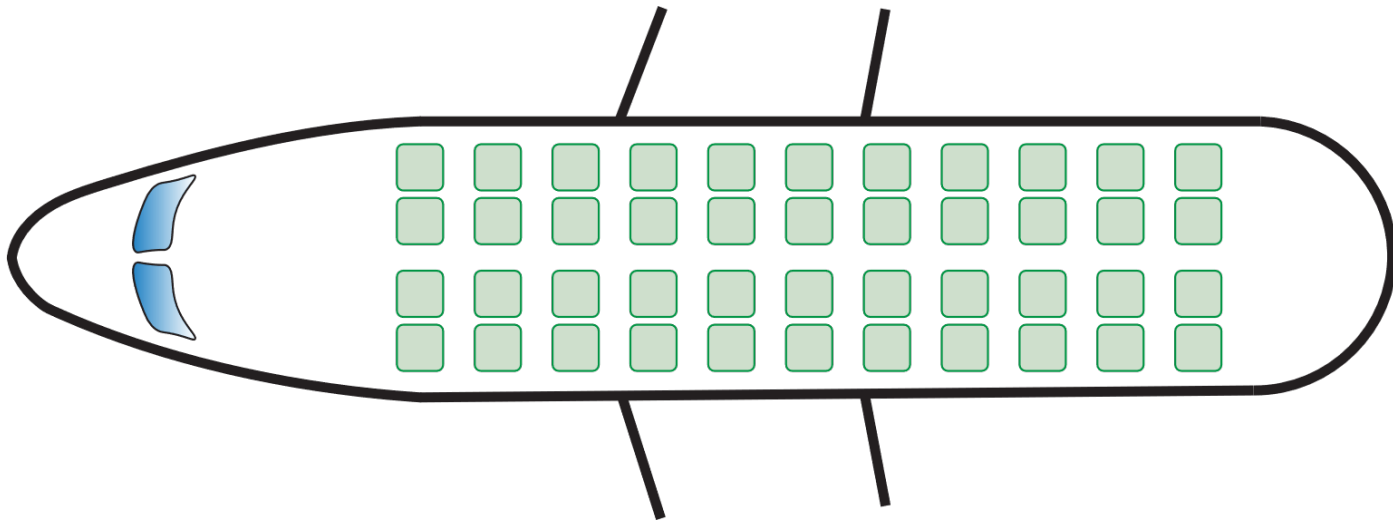
Nonzero



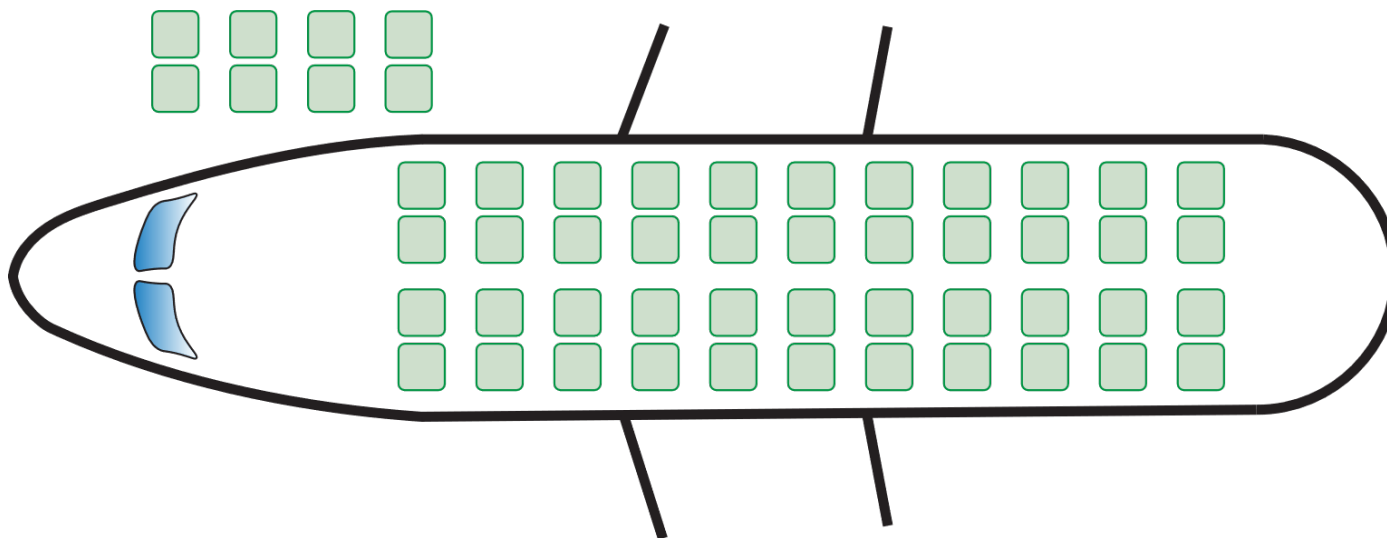
???



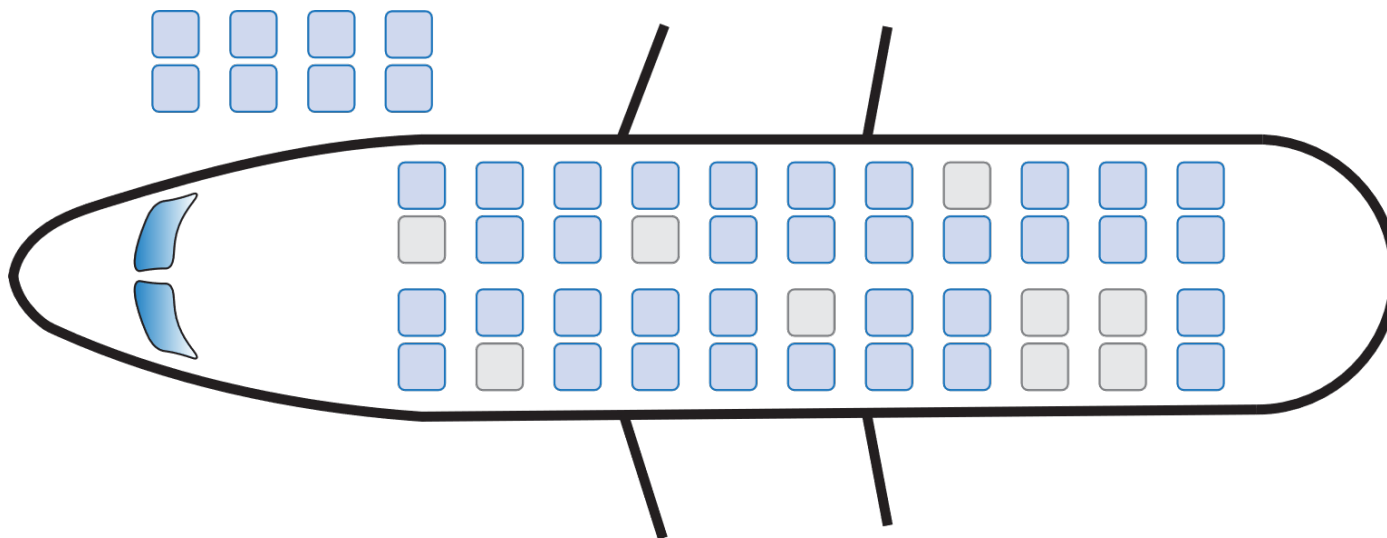
# Overbooking



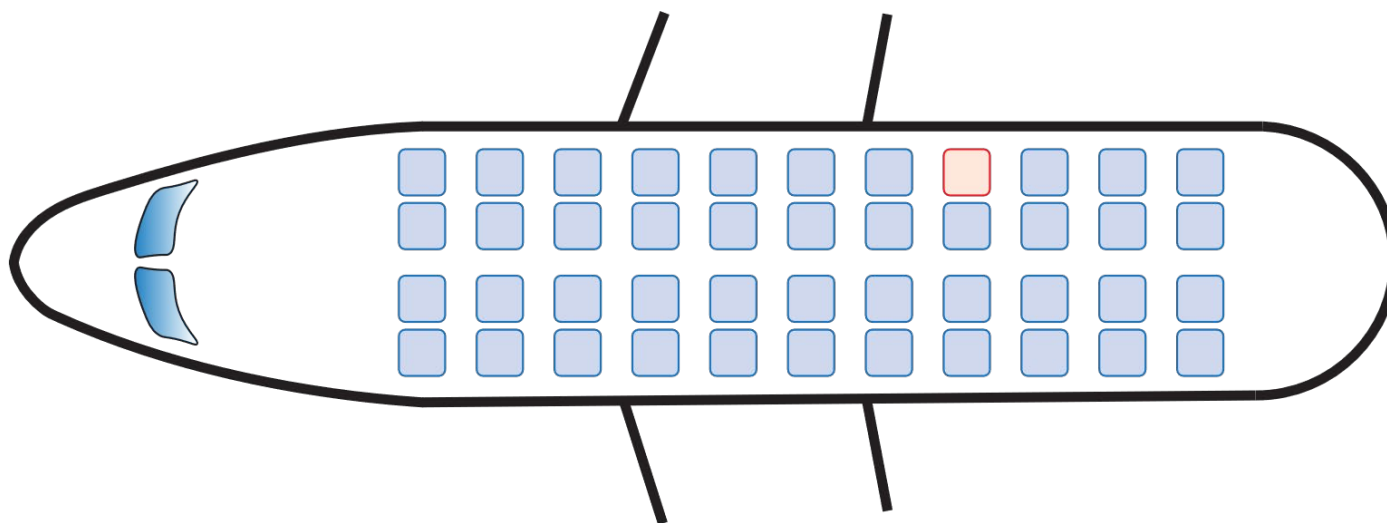
# Overbooking



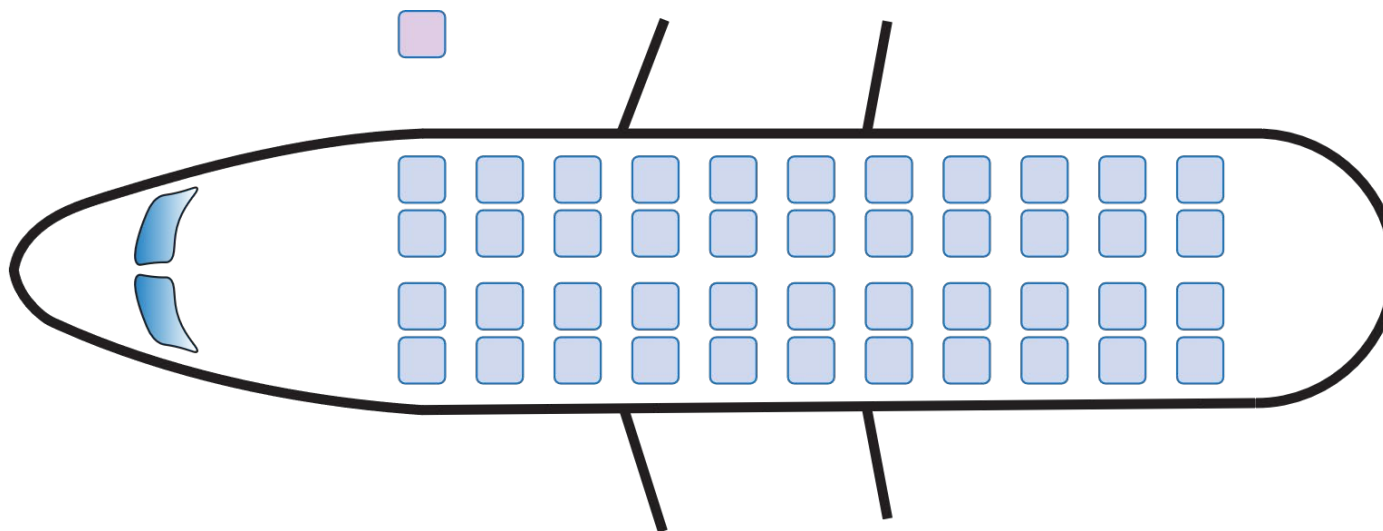
# Overbooking



# Overbooking

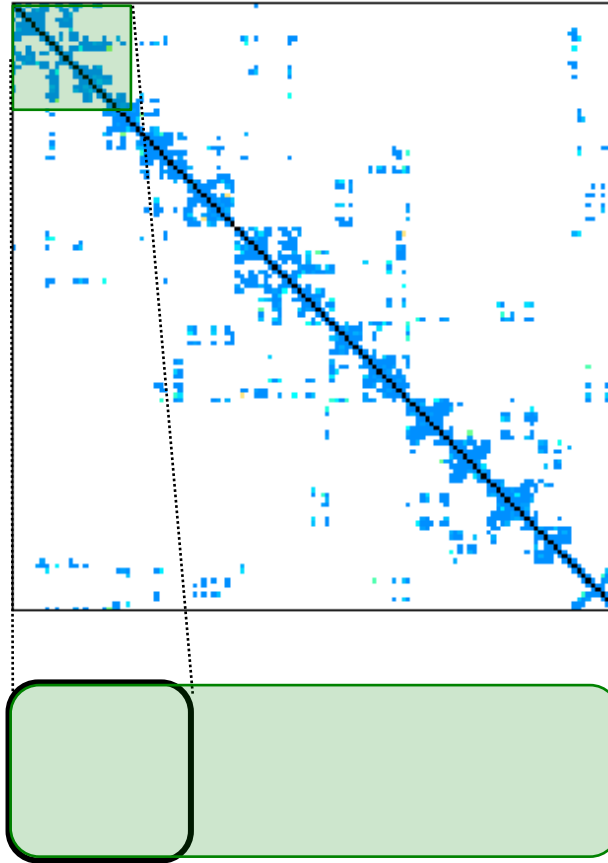


# Overbooking

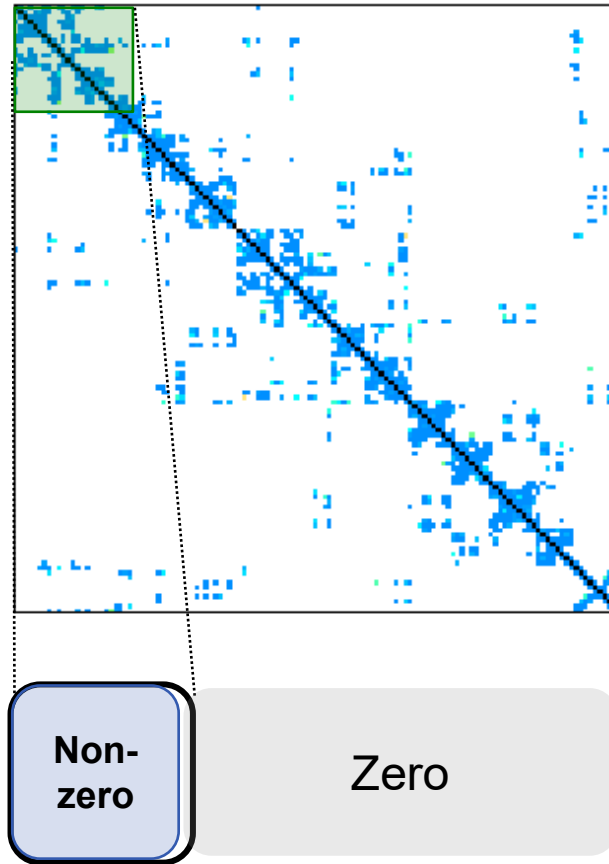




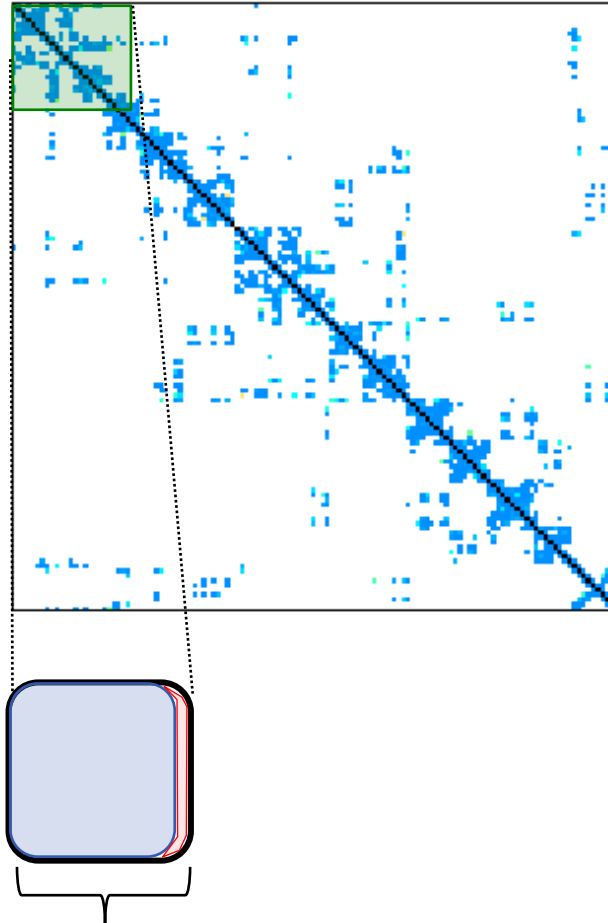
# Overbooking buffers for sparse matrices



# Overbooking buffers for sparse matrices



# Overbooking buffers for sparse matrices



**1300x smaller buffer than prior work**  
**2.4x smaller than optimal coordinate-space tiling**