

# Derin Sinir Ağları için Değişken Boyutlu Sistolik Diziler

Ahmet Caner Yüzügüler – Doktora adayı

Prof. Babak Falsafi - Danışman

# Derin Sinir Ağları (DSA)

- ✓ Hızlı algoritma tasarımı
  - Otomatik öznitelik (feature) çıkarımı
- ✓ İnsan-üstü isabet oranı
  - İmagenet veriseti hata oranı:
- ✗ Çok yüksek bilgisayarım gereksinimi

Konvansiyonel  
görüntü-işleme  
algoritmaları



% 25

vs

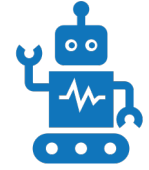
İnsan



% 5

vs

YZ



% 3

*Kaynak: Sze, 2017*

Histogram of oriented  
gradients (HOG):

0.019 GOPs

GooleNet v1:

1.43 GOPs

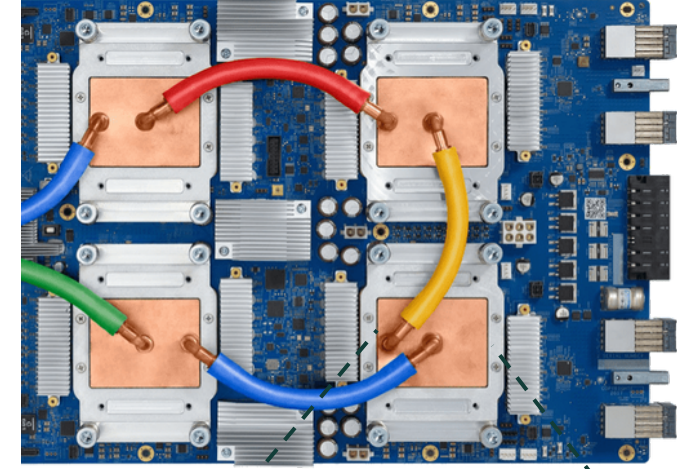
x74

*Kaynak: Suleiman, 2017*

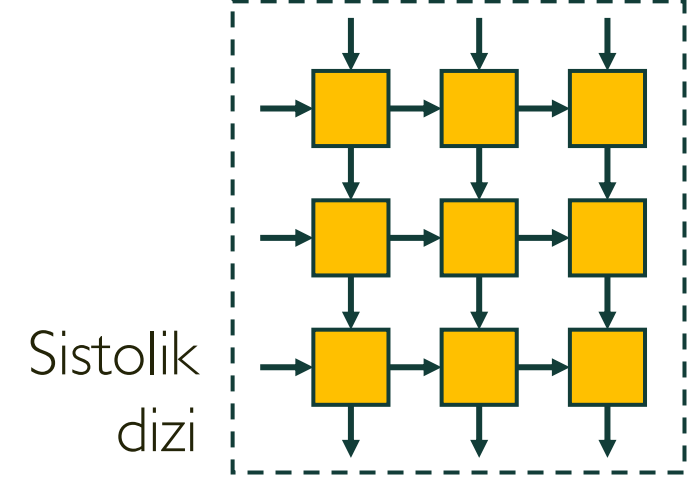
# Sistolik Diziler

- Google'ın DSA'lar için ASIC çözümü:
  - Tensor Processing Unit (TPU)
  - 256 x 256 sistolik dizi
- Neden sistolik dizi?
  - ✓ Yüksek verimli işlem kapasitesi
    - 92 TeraOps/s @ 40 W
  - ✓ Ölçeklenebilirlik
  - ✓ Kolay tasarım – hızlı geliştirme süresi
- Dezavantajı
  - ✗ Veri paylaşımında düşük esneklik
  - ✗ Düşük kapasite kullanım oranı
    - < %30 (GoogleNet)

Resim: <https://cloud.google.com/tpu/>



Google TPU



# Önerdiğimiz Çözüm

- İri parçalı yeniden yapılandırılabilir mimari (Coarse-grained reconfigurable architecture)

Bloklar arasında programlanabilir bağlantılar



Esnek veri aktarım altyapısı



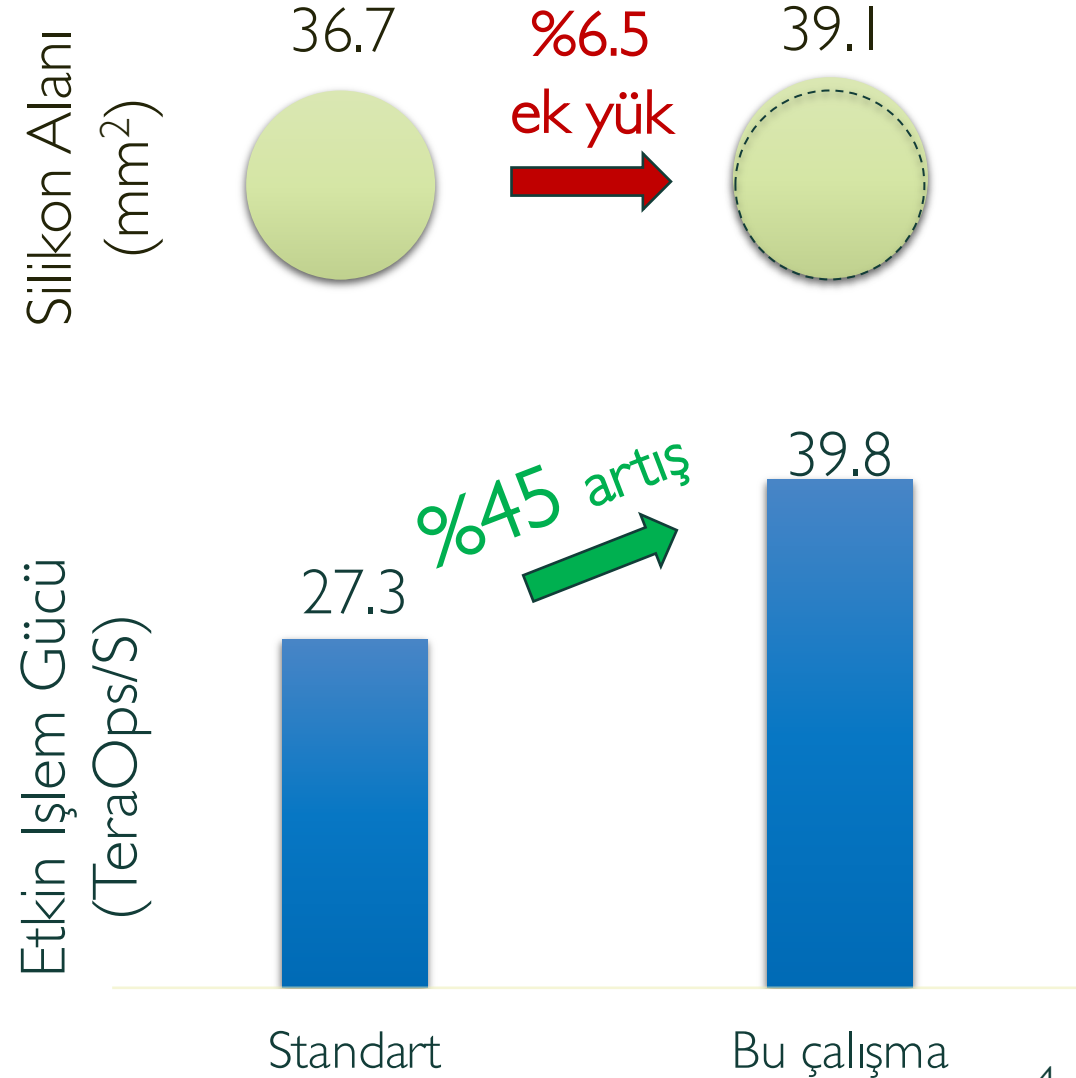
Farklı boyutlardaki DSA katmanlarına daha iyi uyum



Yüksek kapasite kullanımı



Yüksek işlem gücü



# İçerik

- Giriş
- Sistolik diziler
- Düşük kapasite kullanım oranı sebepleri
- Önerilen çözüm
- Sonuçlar

# Sistolik Diziler – Çalışma Prensipleri

- Temelde gerçekleştirdikleri işlem: **matris-matris çarpımı**
- DSA'larını oluşturan işlemlerin büyük çoğunluğu (>%95):
  - **Aktivasyon** ve **filtre** matrislerinin çarpımı

- Örnek:

Aktivasyon matrisi

|          |          |          |
|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

×

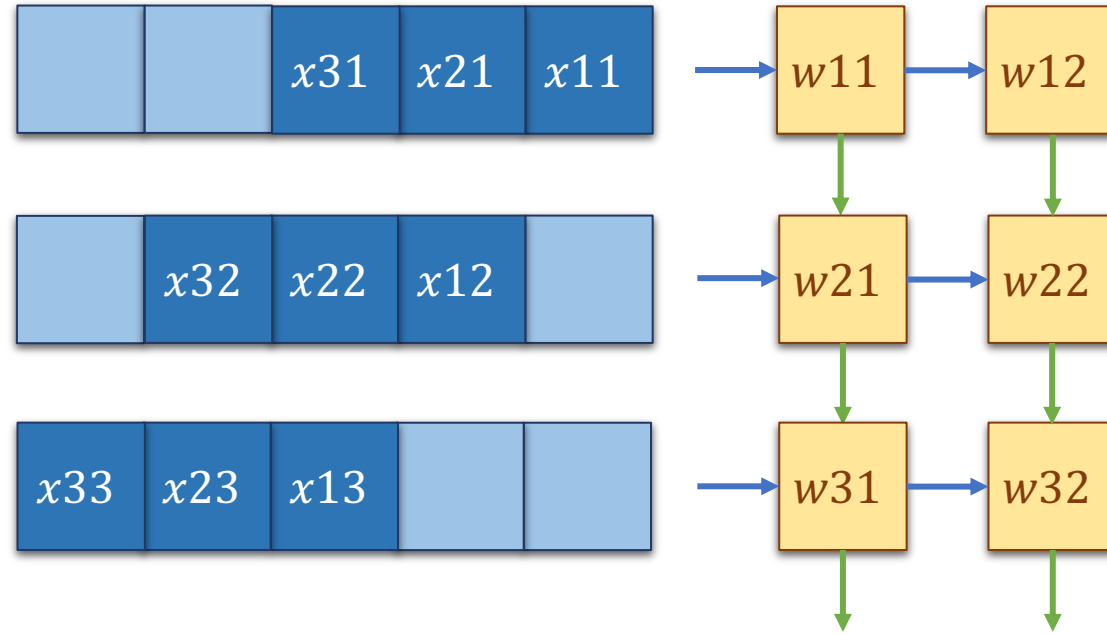
Filtre matrisi

|          |          |
|----------|----------|
| $w_{11}$ | $w_{12}$ |
| $w_{21}$ | $w_{22}$ |
| $w_{31}$ | $w_{32}$ |

=

|          |          |
|----------|----------|
| $y_{11}$ | $y_{12}$ |
| $y_{21}$ | $y_{22}$ |
| $y_{31}$ | $y_{32}$ |

# Sistolik Diziler – Çalışma Prensipleri



Aktivasyon  
matrisi

|                 |                 |                 |
|-----------------|-----------------|-----------------|
| x <sub>11</sub> | x <sub>12</sub> | x <sub>13</sub> |
| x <sub>21</sub> | x <sub>22</sub> | x <sub>23</sub> |
| x <sub>31</sub> | x <sub>32</sub> | x <sub>33</sub> |

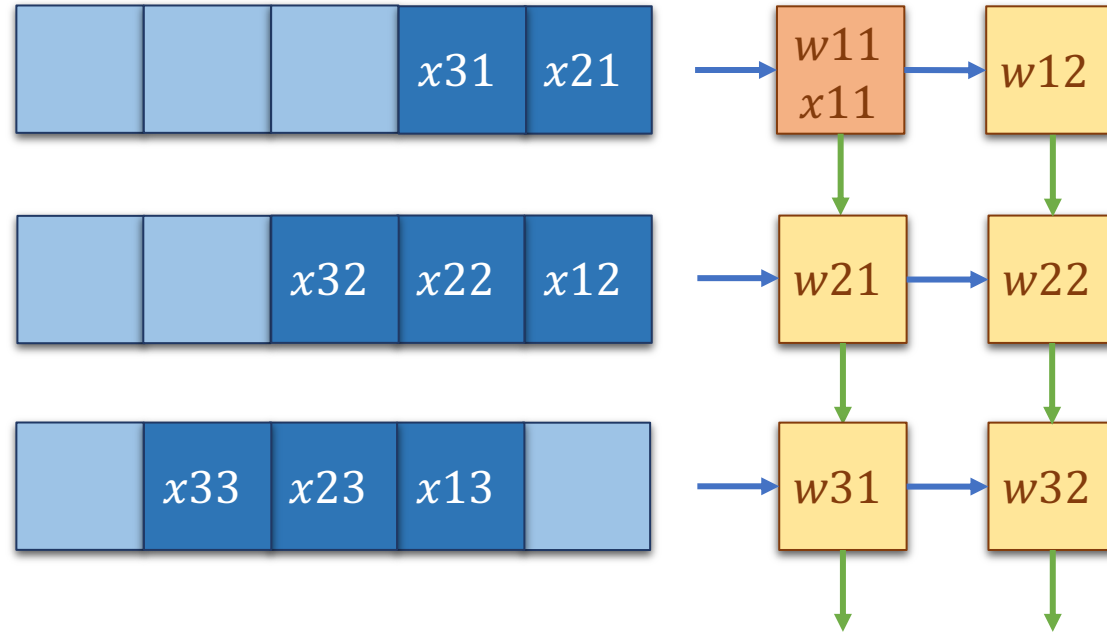
×

Filtre  
matrisi

|                 |                 |
|-----------------|-----------------|
| w <sub>11</sub> | w <sub>12</sub> |
| w <sub>21</sub> | w <sub>22</sub> |
| w <sub>31</sub> | w <sub>32</sub> |

=

# Sistolik Diziler – Çalışma Prensipleri



Aktivasyon  
matrisi

|          |          |          |
|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

×

Filtre  
matrisi

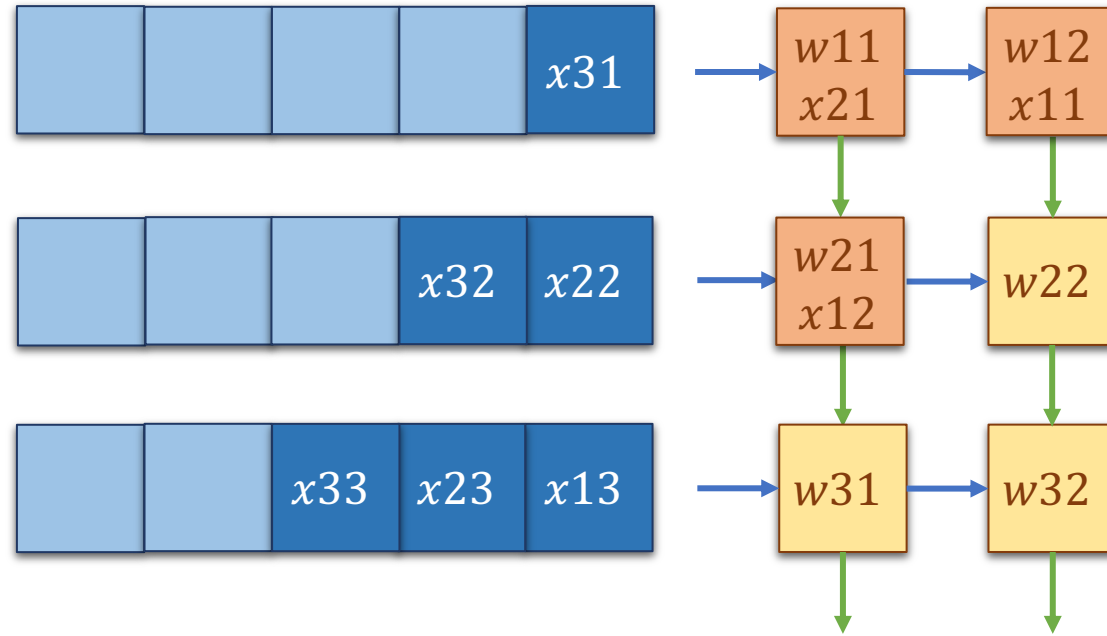
|          |          |
|----------|----------|
| $w_{11}$ | $w_{12}$ |
| $w_{21}$ | $w_{22}$ |
| $w_{31}$ | $w_{32}$ |

=

$$y_{11} = w_{11} x_{11}$$



# Sistolik Diziler – Çalışma Prensipleri



Aktivasyon  
matrisi

|          |          |          |
|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

×

Filtre  
matrisi

|          |          |
|----------|----------|
| $w_{11}$ | $w_{12}$ |
| $w_{21}$ | $w_{22}$ |
| $w_{31}$ | $w_{32}$ |

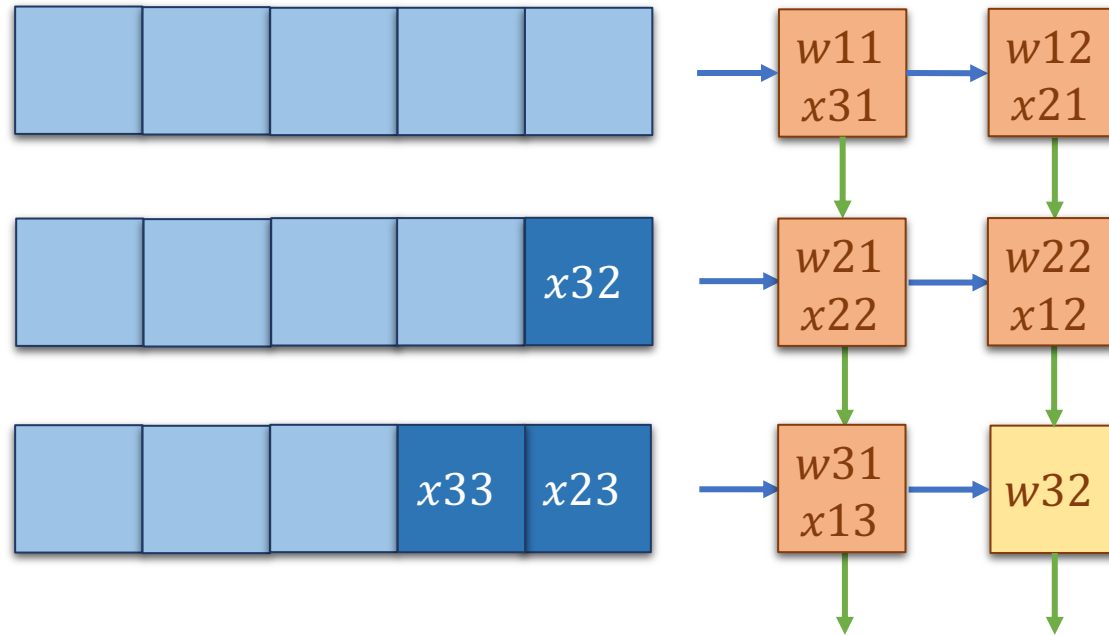
=

$$y_{11} = w_{11} x_{11} + w_{21} x_{12}$$

$$y_{12} = w_{12} x_{11}$$

$$y_{21} = w_{11} x_{21}$$

# Sistolik Diziler – Çalışma Prensipleri



Aktivasyon  
matrisi

|          |          |          |
|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

×

Filtre  
matrisi

|          |          |
|----------|----------|
| $w_{11}$ | $w_{12}$ |
| $w_{21}$ | $w_{22}$ |
| $w_{31}$ | $w_{32}$ |

=

$$y_{11} = w_{11} x_{11} + w_{21} x_{12} + w_{31} x_{13}$$

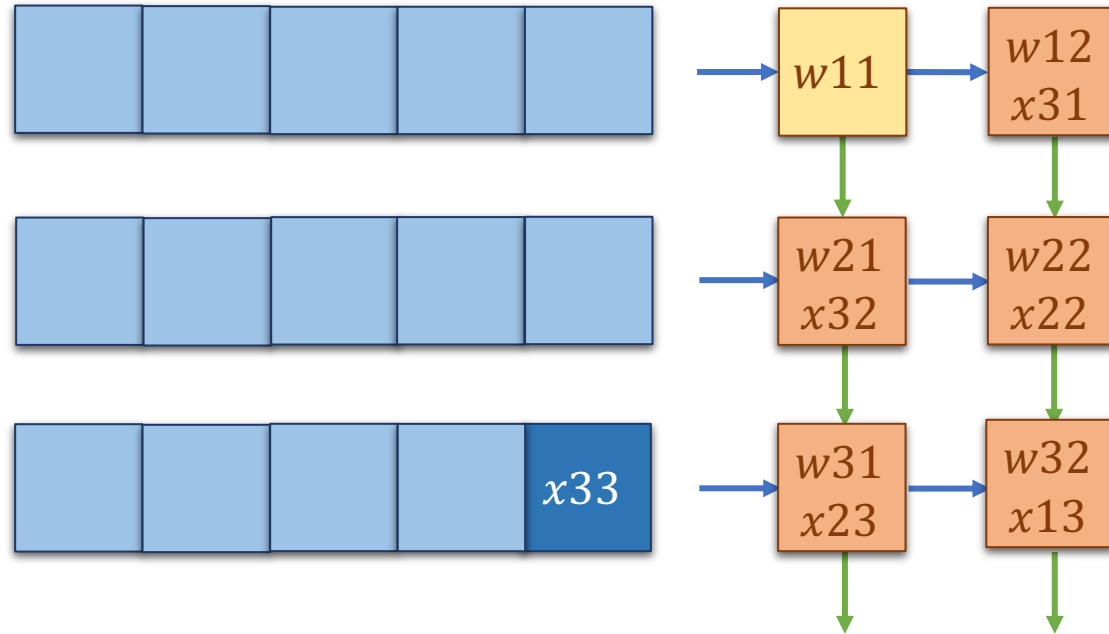
$$y_{12} = w_{12} x_{11} + w_{22} x_{12}$$

$$y_{21} = w_{11} x_{21} + w_{21} x_{22}$$

$$y_{22} = w_{12} x_{21}$$

$$y_{31} = w_{11} x_{31}$$

# Sistolik Diziler – Çalışma Prensipleri



Aktivasyon  
matrisi

|          |          |          |
|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

×

Filtre  
matrisi

|          |          |
|----------|----------|
| $w_{11}$ | $w_{12}$ |
| $w_{21}$ | $w_{22}$ |
| $w_{31}$ | $w_{32}$ |

=

$$y_{11} = w_{11} x_{11} + w_{21} x_{12} + w_{31} x_{13}$$

$$y_{12} = w_{12} x_{11} + w_{22} x_{12} + w_{32} x_{13}$$

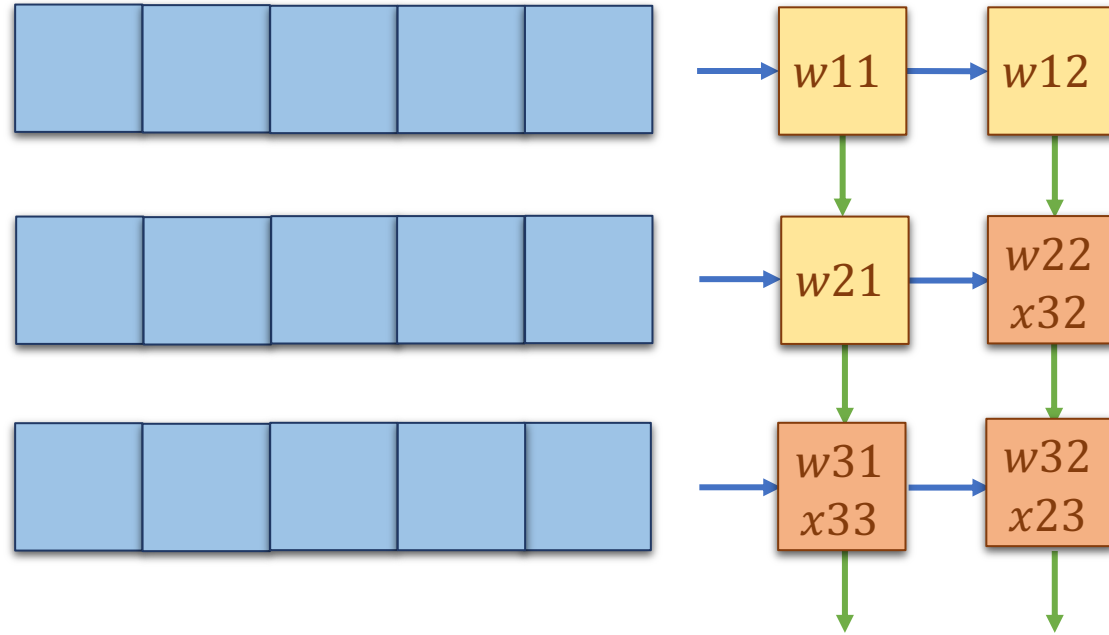
$$y_{21} = w_{11} x_{21} + w_{21} x_{22} + w_{31} x_{23}$$

$$y_{22} = w_{12} x_{21} + w_{22} x_{22}$$

$$y_{31} = w_{11} x_{31} + w_{21} x_{32}$$

$$y_{32} = w_{12} x_{31}$$

# Sistolik Diziler – Çalışma Prensipleri



Aktivasyon  
matrisi

|          |          |          |
|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

×

Filtre  
matrisi

|          |          |
|----------|----------|
| $w_{11}$ | $w_{12}$ |
| $w_{21}$ | $w_{22}$ |
| $w_{31}$ | $w_{32}$ |

=

$$y_{11} = w_{11} x_{11} + w_{21} x_{12} + w_{31} x_{13}$$

$$y_{12} = w_{12} x_{11} + w_{22} x_{12} + w_{32} x_{13}$$

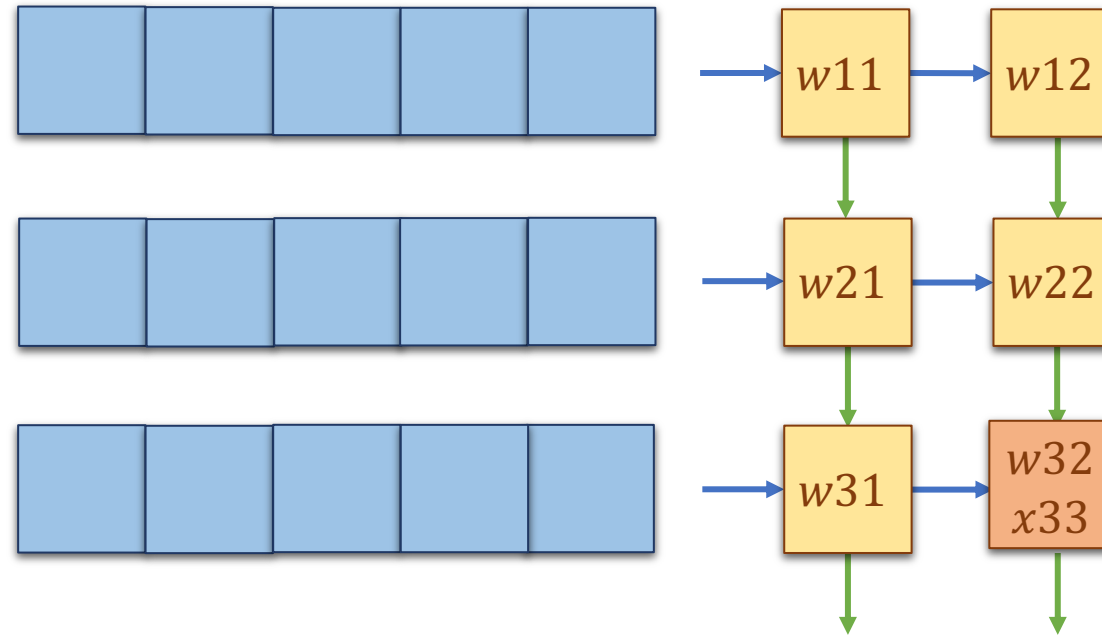
$$y_{21} = w_{11} x_{21} + w_{21} x_{22} + w_{31} x_{23}$$

$$y_{22} = w_{12} x_{21} + w_{22} x_{22} + w_{32} x_{23}$$

$$y_{31} = w_{11} x_{31} + w_{21} x_{32} + w_{31} x_{33}$$

$$y_{32} = w_{12} x_{31} + w_{22} x_{32}$$

# Sistolik Diziler – Çalışma Prensipleri



Aktivasyon  
matrisi

|          |          |          |
|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

×

Filtre  
matrisi

|          |          |
|----------|----------|
| $w_{11}$ | $w_{12}$ |
| $w_{21}$ | $w_{22}$ |
| $w_{31}$ | $w_{32}$ |

=

$$y_{11} = w_{11} x_{11} + w_{21} x_{12} + w_{31} x_{13}$$

$$y_{12} = w_{12} x_{11} + w_{22} x_{12} + w_{32} x_{13}$$

$$y_{21} = w_{11} x_{21} + w_{21} x_{22} + w_{31} x_{23}$$

$$y_{22} = w_{12} x_{21} + w_{22} x_{22} + w_{32} x_{23}$$

$$y_{31} = w_{11} x_{31} + w_{21} x_{32} + w_{31} x_{33}$$

$$y_{32} = w_{12} x_{31} + w_{22} x_{32} + w_{32} x_{33}$$

# Sistolik Dizilerin Avantaj/Dezavantajları

✓ Veri aktarımı: komşu işlemci birimleri arasında

- Veri aktarımı için harcanan enerji  $\rightarrow 0$
- Kısa bağlantı hatları  $\rightarrow$  Ölçeklenebilir tasarım

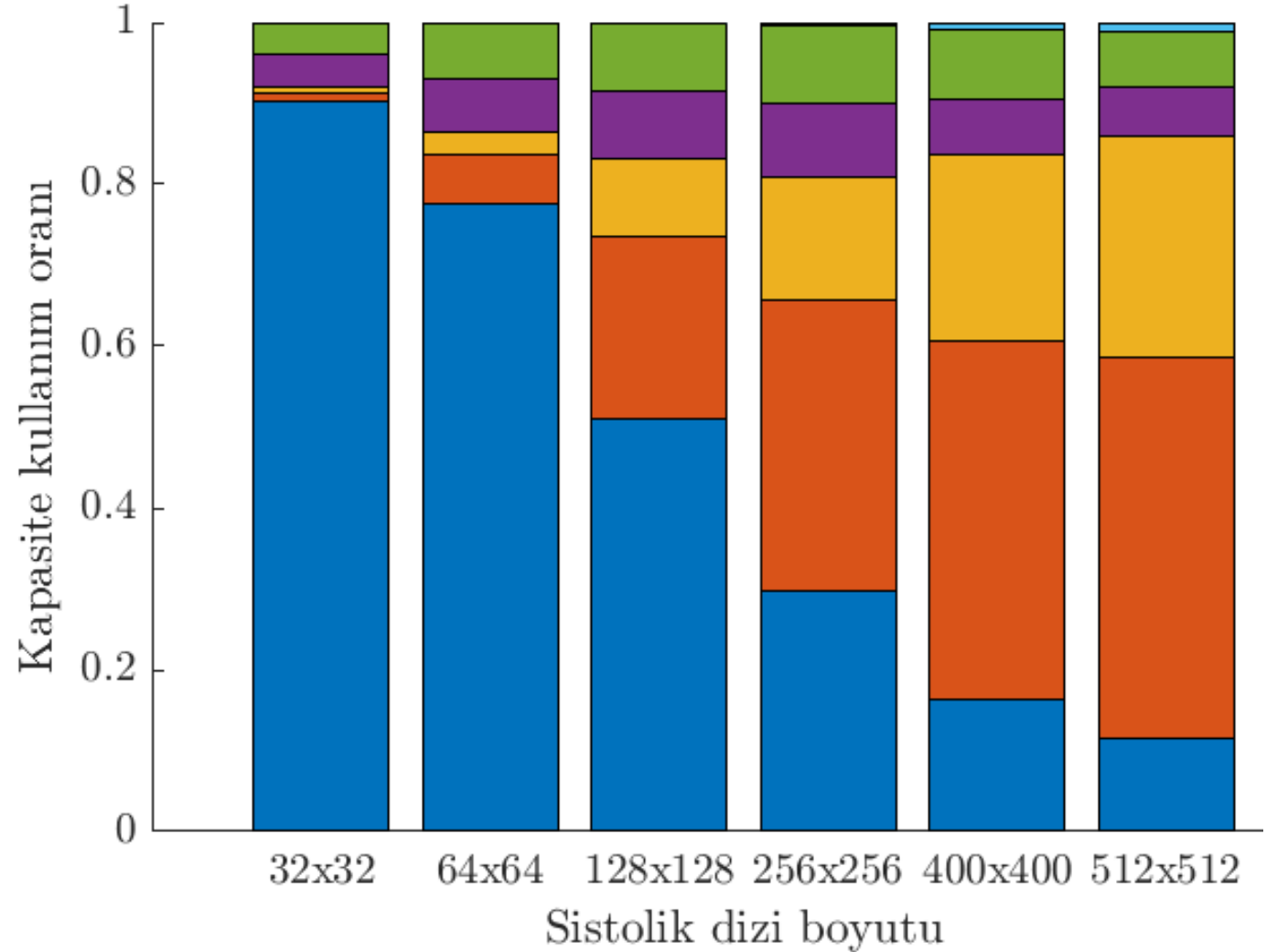
✗ Düşük kapasite kullanım oranı

1. Yatay aktarım gecikmesi
2. Dikey aktarım gecikmesi
3. Dar filtre boşluğu
4. Sığ filtre boşluğu

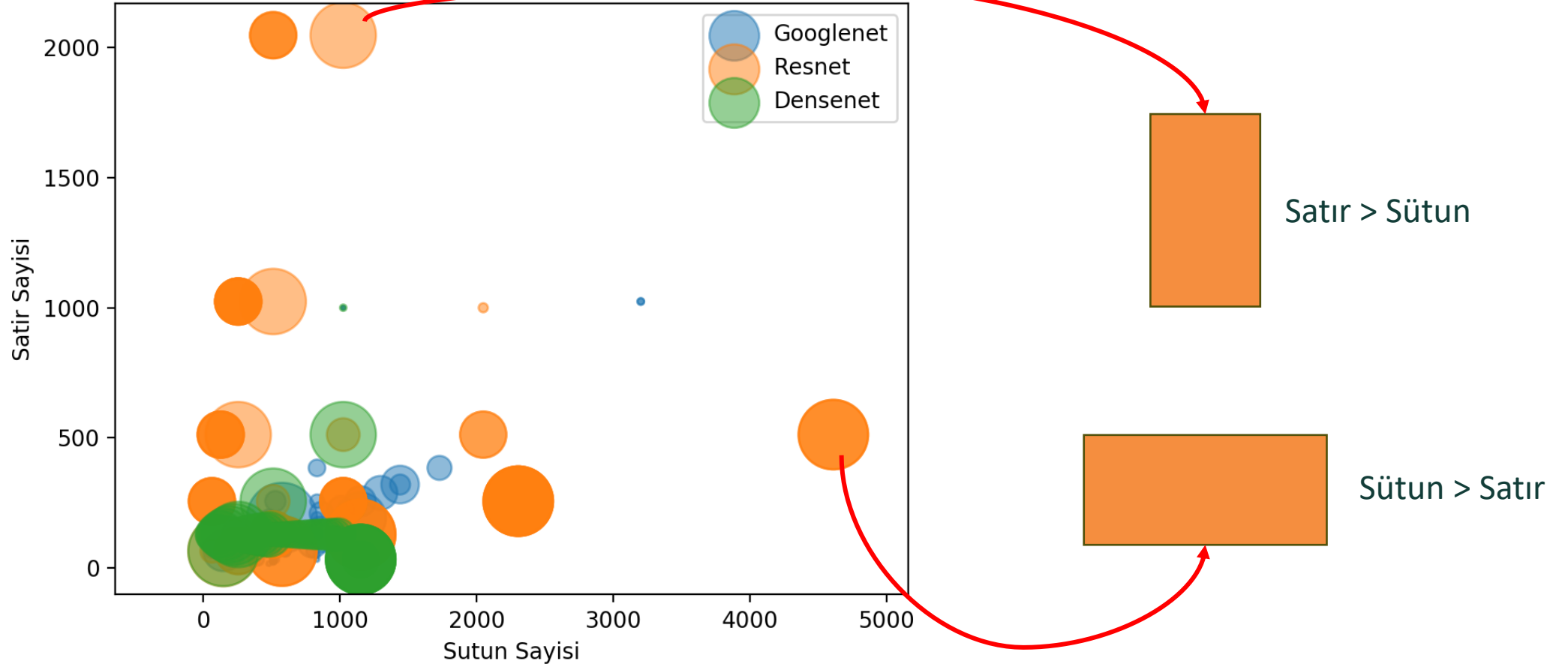


# Düşük Kapasite Kullanımı Nedenleri

- Büyük sistolik dizi → Düşük kapasite kullanımı
- Kapasitesin **>%50**'sinin kullanılmama sebebi:
  - Dar/sığ filtre boşluğu



# DSA Katman Boyutlarındaki Değişkenlik

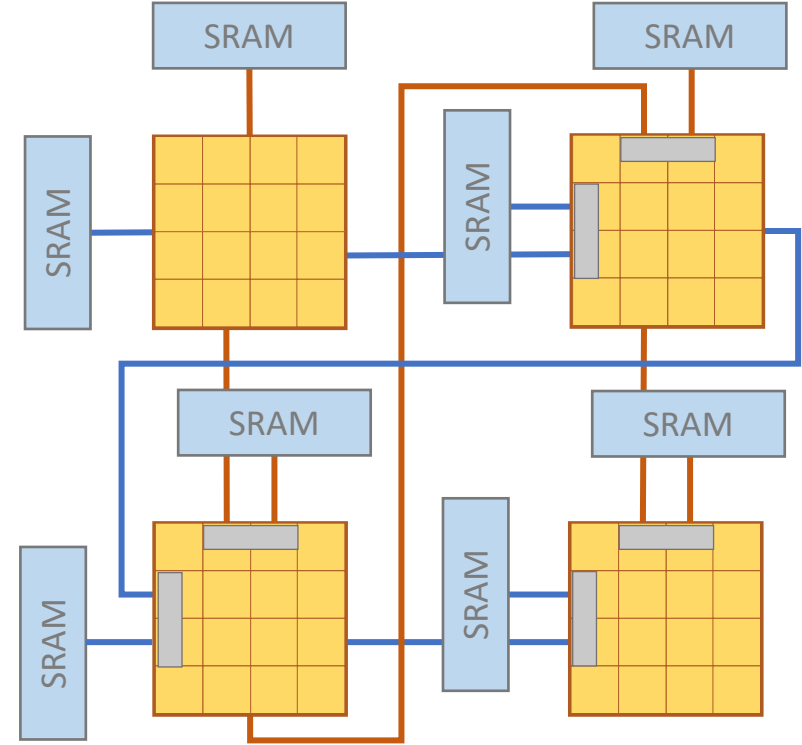


Sabit sistolik dizi + değişken filtre boyutu = düşük kapasite kullanımı

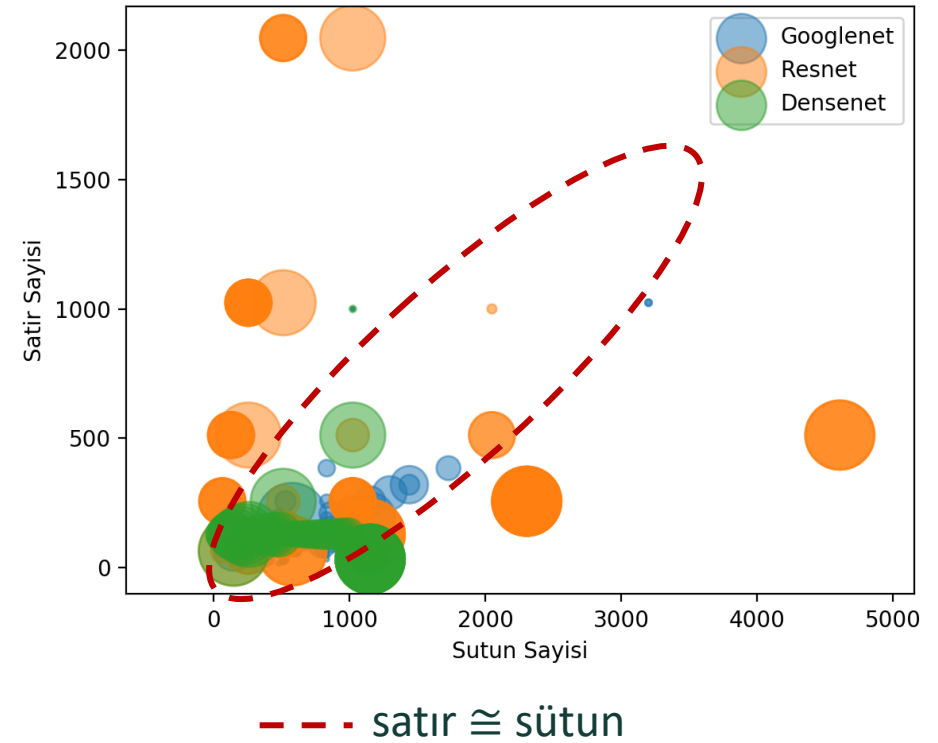
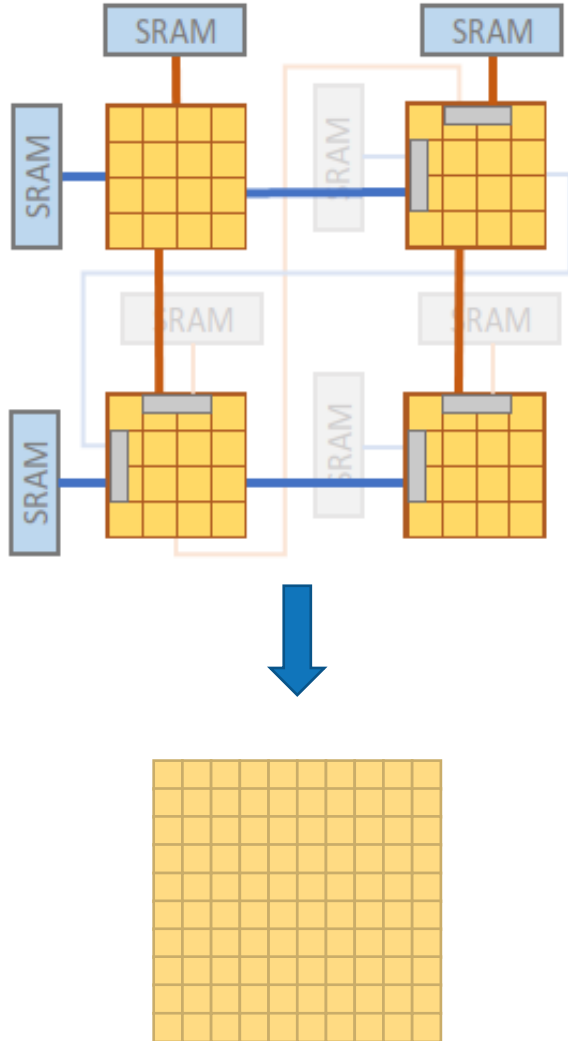


# Önerdiğimiz Çözüm

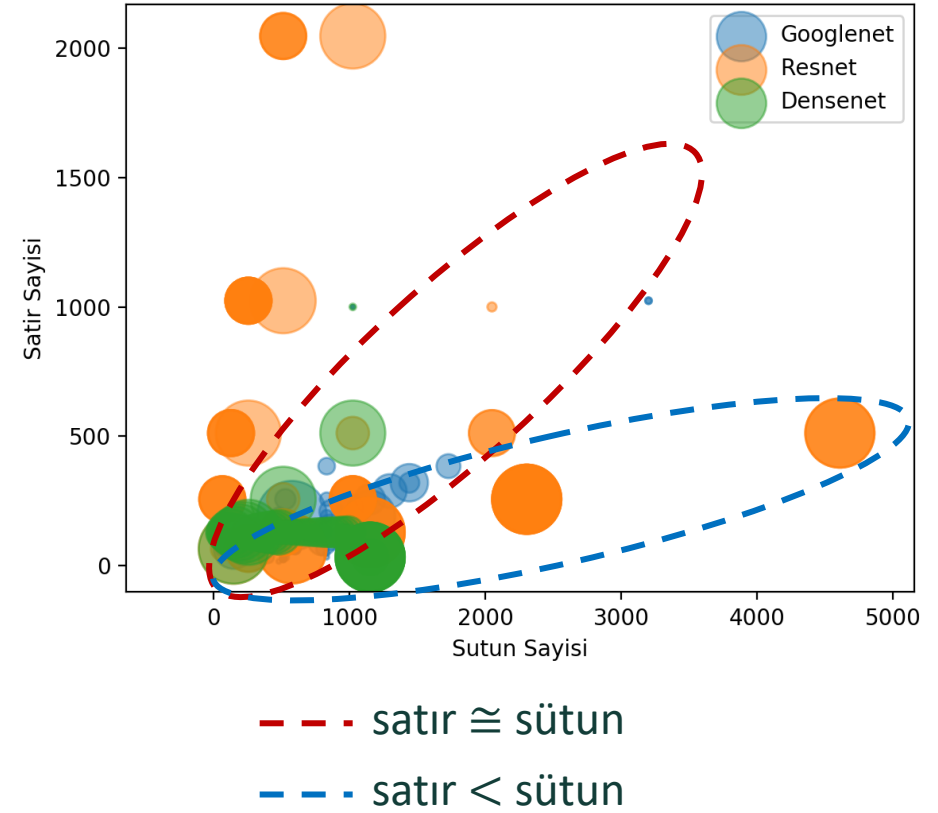
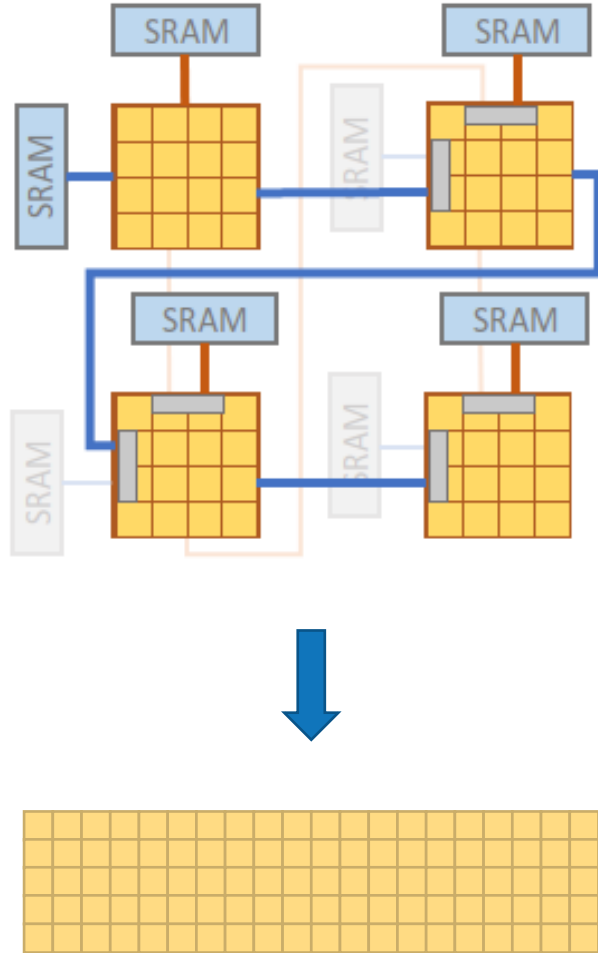
- İri parçalı yeniden yapılandırılabilir mimari (Coarse-grained reconfigurable arch.)
  - Tek büyük parça → Birçok küçük parça
  - Parçalar arası programlanabilir bağlantılar
  - Farklı bağlantı konfigürasyonları → Değişken sistolik dizi boyutu
  - Değişken boyut → Farklı filtrelere daha iyi uyum



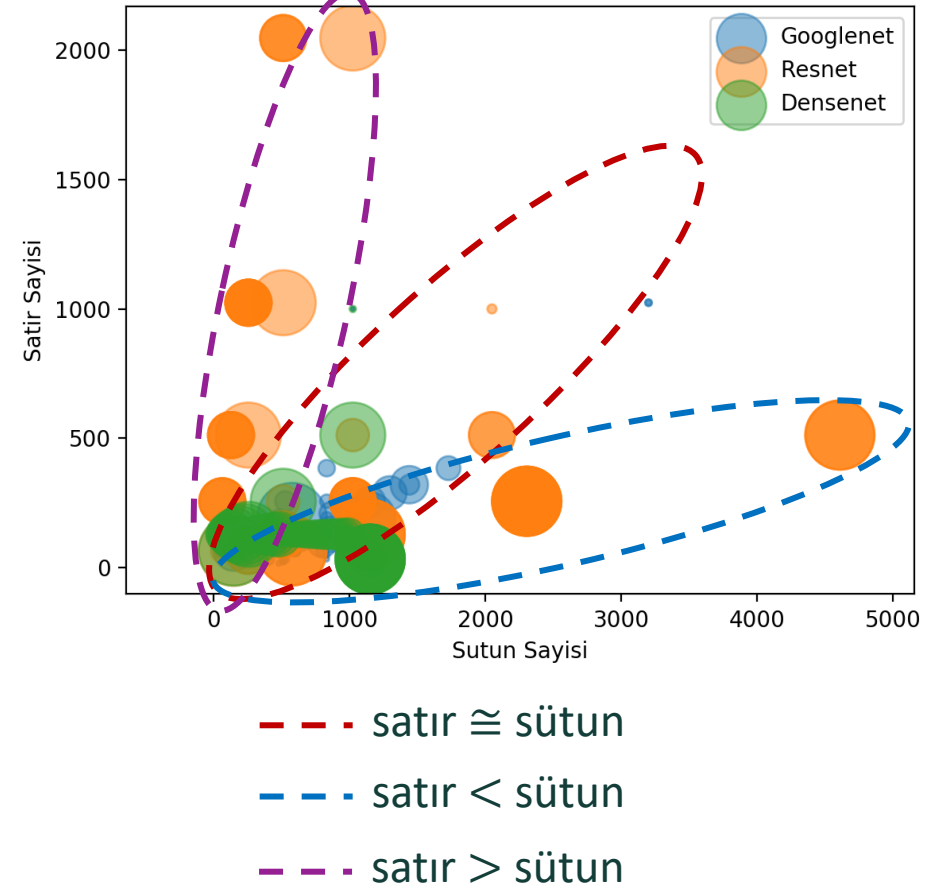
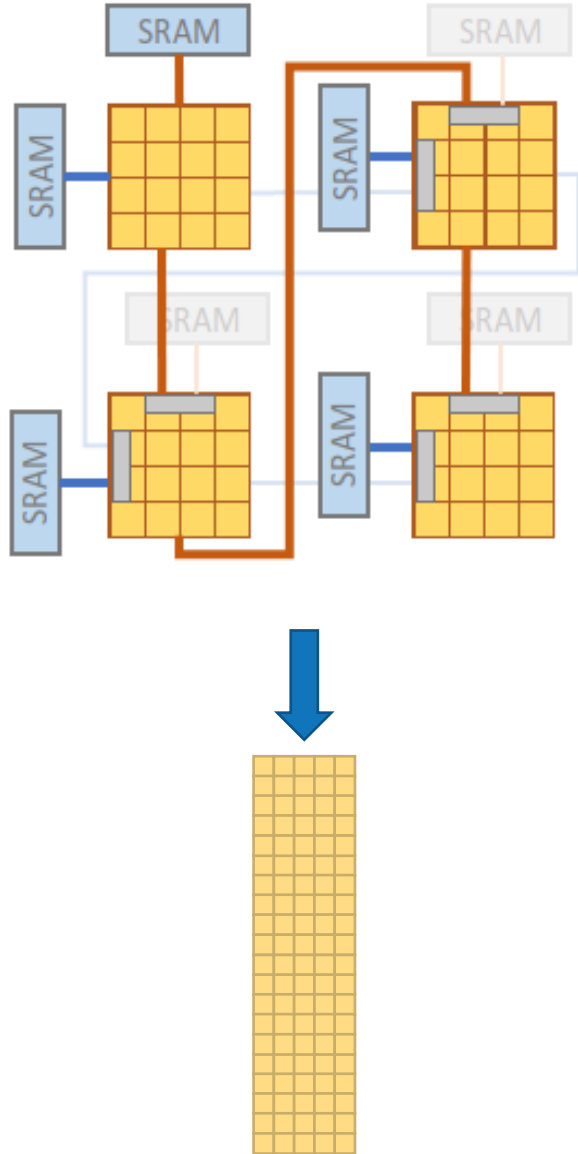
# Kare konfigürasyon



# Geniş konfigürasyon

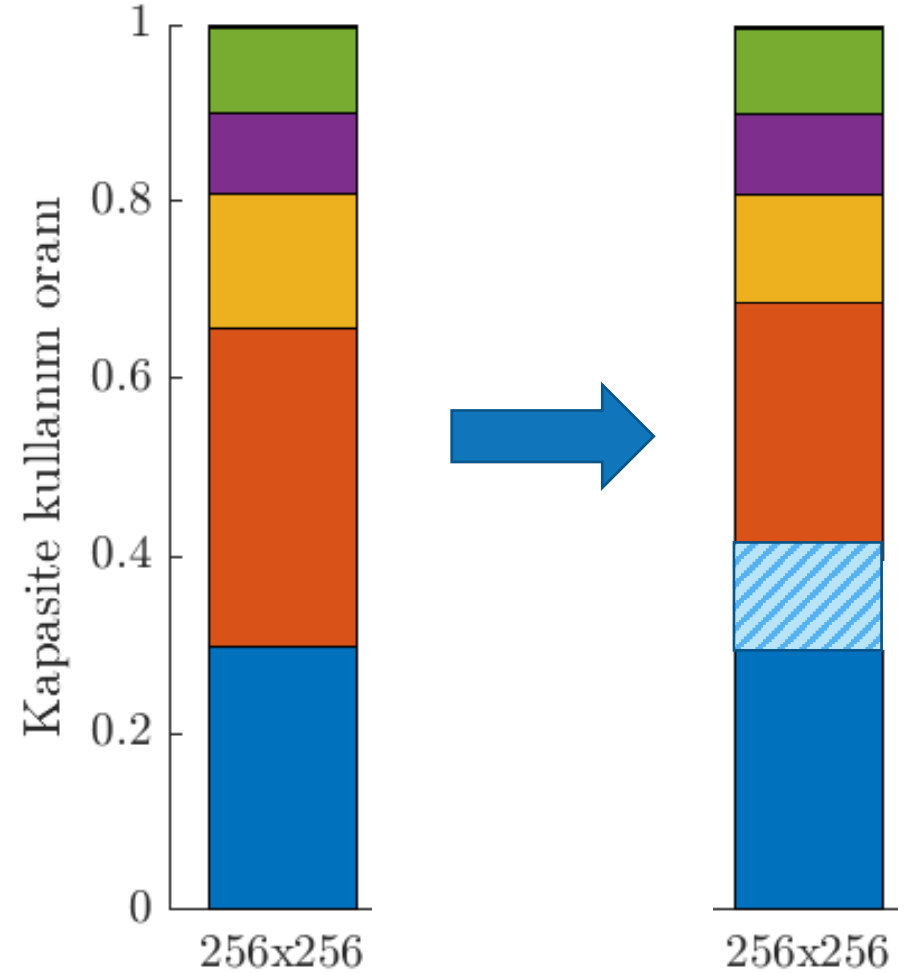


# Derin konfigürasyon



# Sonuçlar

- 8 parçaya bölünmüş 256x256'lik bir sistolik dizide
  - Farklı filtre boyutlarına **%20** daha iyi uyum
  - **%45** daha yüksek etkin işlem gücü



## ■ Google TPU

- 256x256 standart sistolik dizi
- 28nm TSMC
- 36.7 mm<sup>2</sup> silikon alanı

**%6.5 ek yük**

## ■ Bu çalışma

- 8 parçalı 256x256 sistolik dizi
- 28nm TSMC
- 39.1 mm<sup>2</sup> silikon alanı
  - İşlemciler arası bağlantı karmaşıklığı
  - Programlanabilirlik için gerekli çoklayıcılar (multiplexer)
- %43.3 kapasite kullanım oranı
- 39.8 TeraOps/s etkin işlem gücü

- %29.7 kapasite kullanım oranı
- 27.3 TeraOps/s etkin işlem gücü

**%45 artış**

- **Sze,2017:** V. Sze et. al., Efficient Processing of Deep Neural Networks: A Tutorial and Survey, Proceedings of the IEEE, Volume: 105 Issue: 12, 2017
- **Suleiman,2017:** A. Suleiman et. al., Towards Closing the Energy Gap Between HOG and CNN Features for Embedded Vision, IEEE International Symposium on Circuits and Systems (ISCAS), 2017
- N. P. Jouppi et. al., In-Datcenter Performance Analysis of a Tensor Processing Unit, International Symposium on Computer Architecture (ISCA), 2017
- H. T. Kung and C. E. Leiserson, "Systolic Arrays (For VLSI)", Sparse Matrix Proceedings, 1978
- C. Szegedy et. al. Going Deeper with Convolutions, IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015.