



IEEE ICCABS

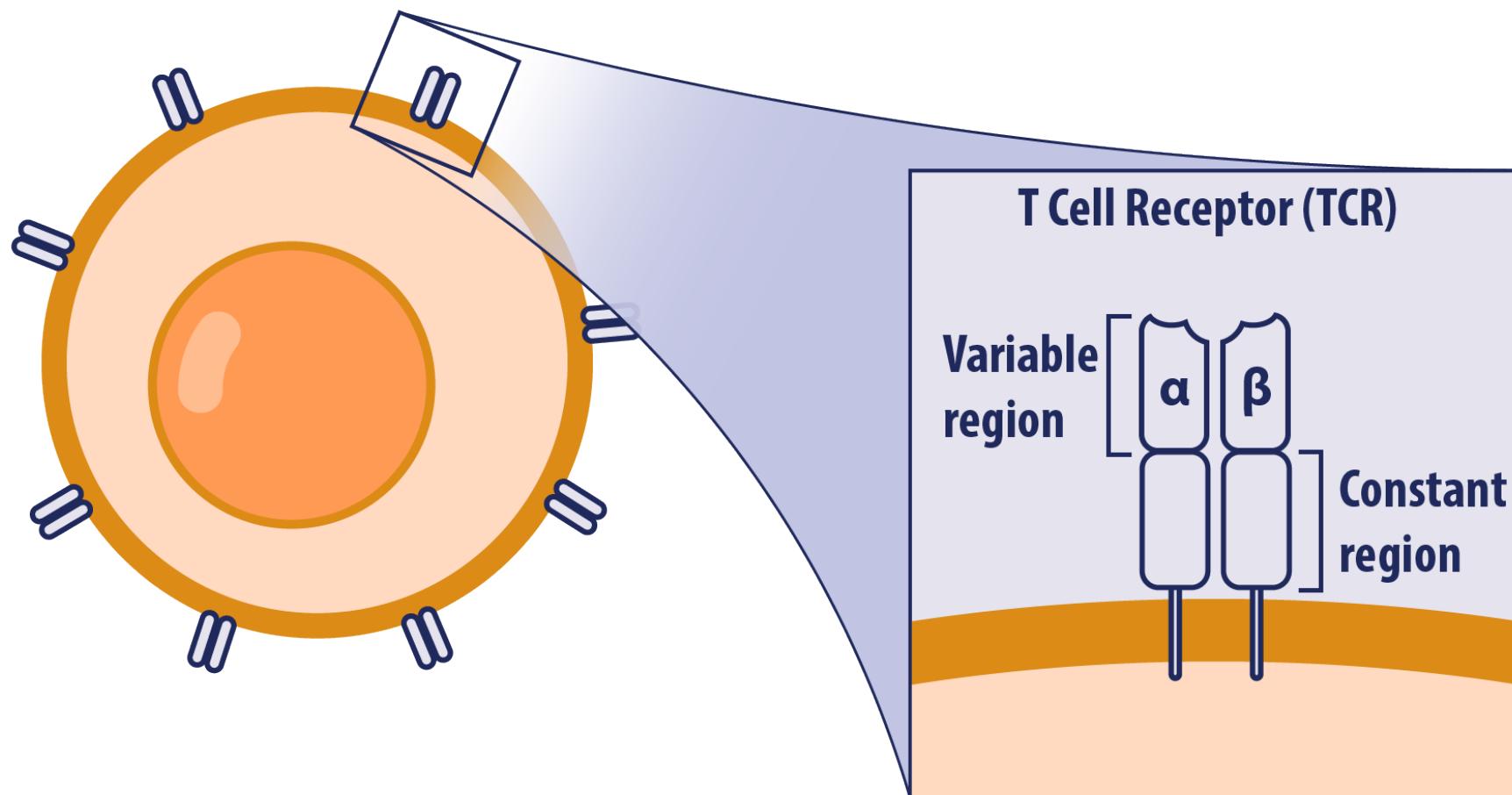
UCONN  
UNIVERSITY OF CONNECTICUT

# Identifying $\alpha\beta$ T cell clones via pooling and b-matching\*

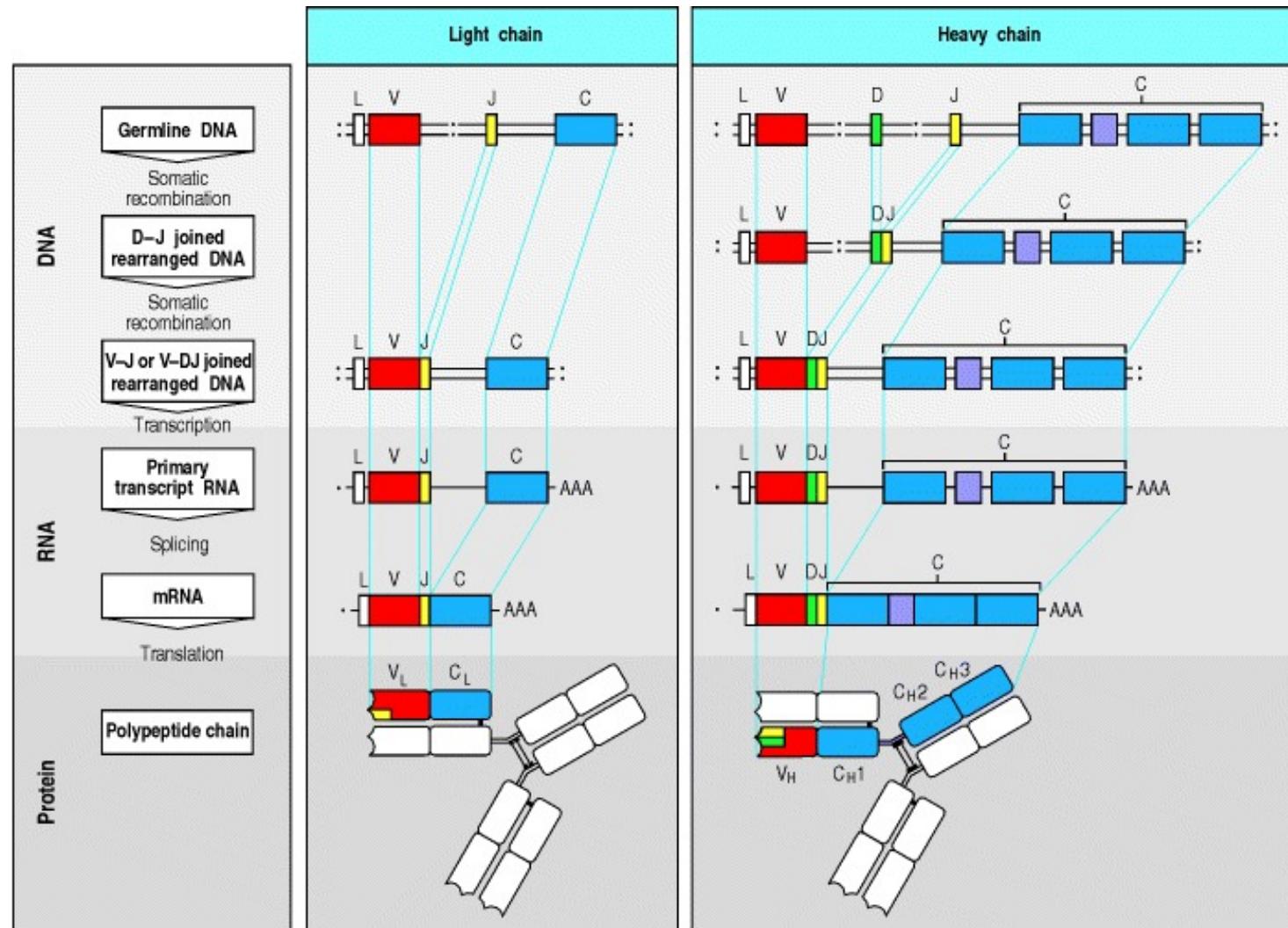
Tyler Daddio, Ion Măndoiu  
University of Connecticut

ICCABS 2017

# $\alpha\beta$ T cells



# Formation of T cell receptors



# Multiple rearrangements?

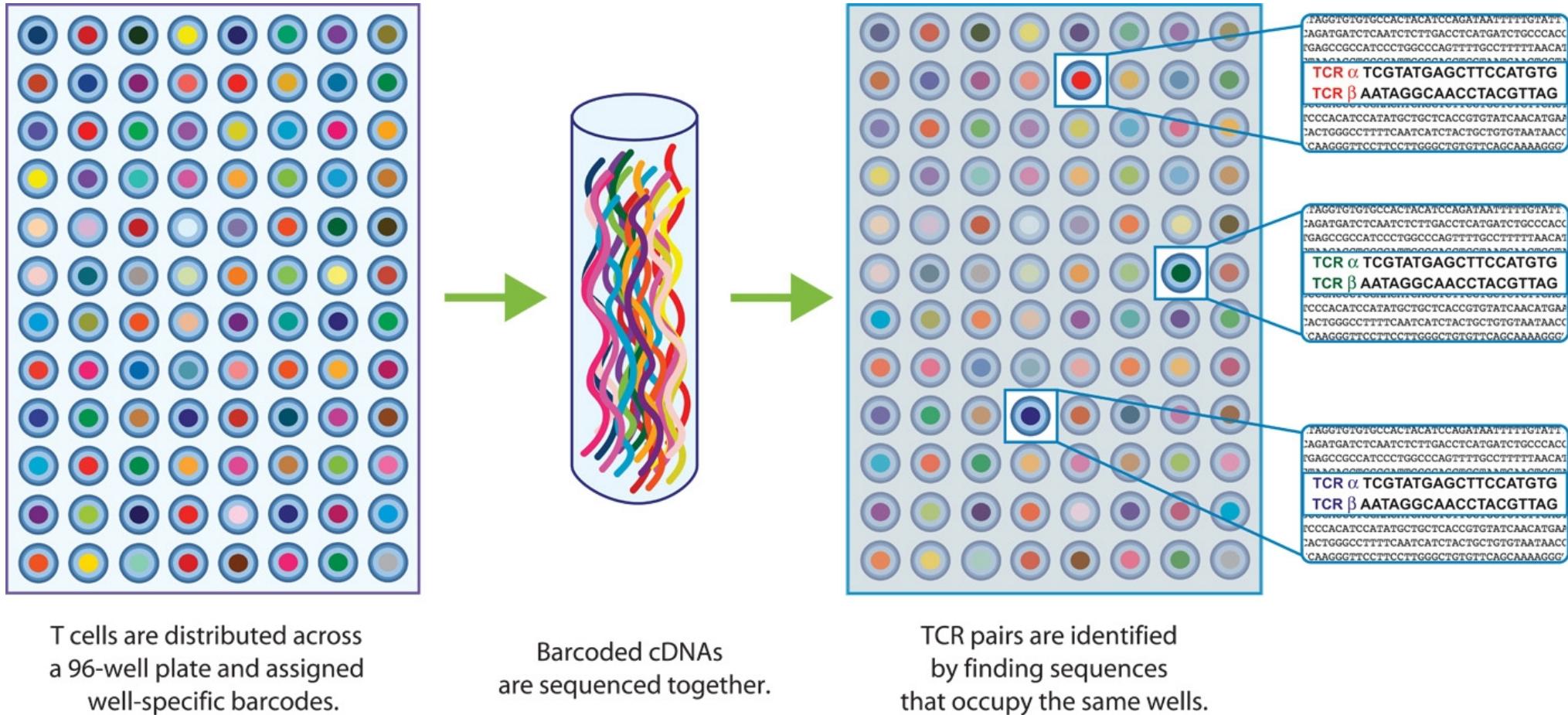
- Up to 30% of T cells contain two in-frame  $\alpha$ -recombinants (i.e. dual alpha)
- Up to 2-10% contain two in-frame  $\beta$ -recombinants (i.e. dual beta)
- Non-productive rearrangements also occur
  - Challenge or opportunity?

# Sharing recombinants?

- Clones can share a single  $\beta$ -recombinant due to:
  - Convergent recombination
  - Maturation process
- Clones can only share a single  $\alpha$ -recombinant due to convergent recombination
- Amount of sharing varies significantly
  - anywhere from 0% to ~50% of  $\alpha$ -recombinants
  - anywhere from 0% to ~40% of  $\beta$ -recombinants



# PairSEQ



# Clone Identification Problem (CIP)

- **Input:**

- observed alpha recombinants  $\alpha_1, \alpha_2, \dots, \alpha_n$
- observed beta recombinants  $\beta_1, \beta_2, \dots, \beta_m$
- recombinant well occurrences (i.e. well subsets)

$$W : \{\alpha_1, \dots, \alpha_n\} \cup \{\beta_1, \dots, \beta_m\} \rightarrow P(\{1, 2, \dots, w\})$$

- **Output:**

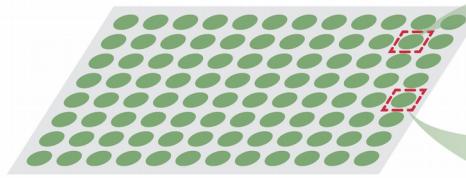
- set of identified T cell clones
  - singles, duals, sharing

# Challenges

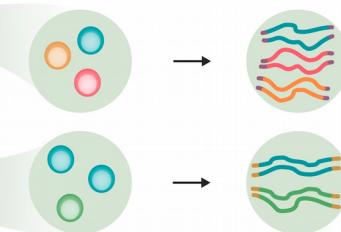
0. Pairing “normal” recombinants? (i.e. *singles*)
1. Multiple rearrangements? (i.e. *duals*)
2. Shared recombinants? (i.e. *sharing*)
3. Computational efficiency?
4. Elegance?

# ALPHABETR

**A** T cells are sampled onto 96-well plates at 10-300 cells/well



**B** cDNA libraries are created from each well with RT-PCR



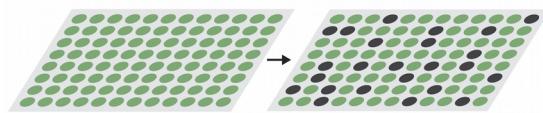
**C** Unpaired CDR3 $\alpha$  and CDR3 $\beta$  are sequenced in each well

$\alpha_1$  CAVTGGDKLIF  
 $\alpha_2$  CALDGGDKIIF  
 $\beta_1$  CASGLARAEQYF  
 $\beta_2$  CASSEGDKVIF

$\alpha_1$  CAVTGGDKLIF  
 $\alpha_3$  CAVTYGYLNF  
 $\alpha_4$  CALTASGLTF  
 $\beta_1$  CASGLARAEQYF  
 $\beta_2$  CASSHSRYEQYF  
 $\beta_3$  CSEVHTARTQYF

**D** A resampling strategy is used to obtain a list of possible TCR pairs by repeatedly performing steps (i), (ii), (iii)

(i) A subset of wells is randomly selected for steps (ii) and (iii)



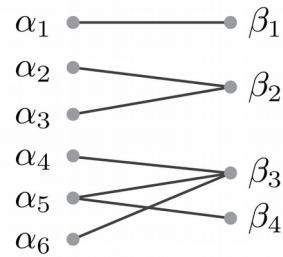
(ii) Association scores are calculated for every  $\alpha$  and  $\beta$  chain across all wells in the subset

$$S_{ij} = \sum_{k=1}^W \left( \frac{\delta_{ij}^k}{c_\alpha^k} + \frac{\delta_{ij}^k}{c_\beta^k} \right)$$

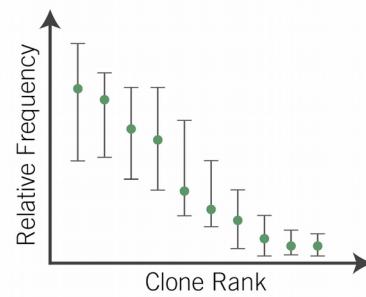
(iii) Scores are used to select likely  $\alpha\beta$  pairs within each well

	$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$
$\beta_1$	24.0	0.6	1.2	4.3	8.2
$\beta_2$	0.4	0.2	60.2	0.7	2.2
$\beta_3$	1.0	0.2	1.2	9.0	3.0
$\beta_4$	3.2	30.1	0.1	0.4	2.1

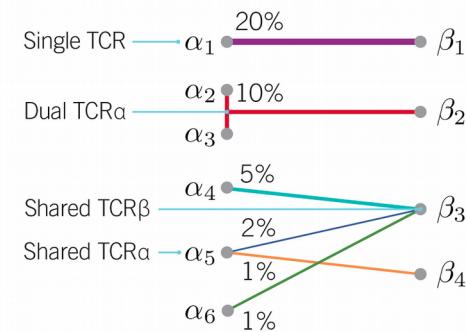
**E** Pairs from Step D present in more than a threshold proportion of replicates are candidate TCRs



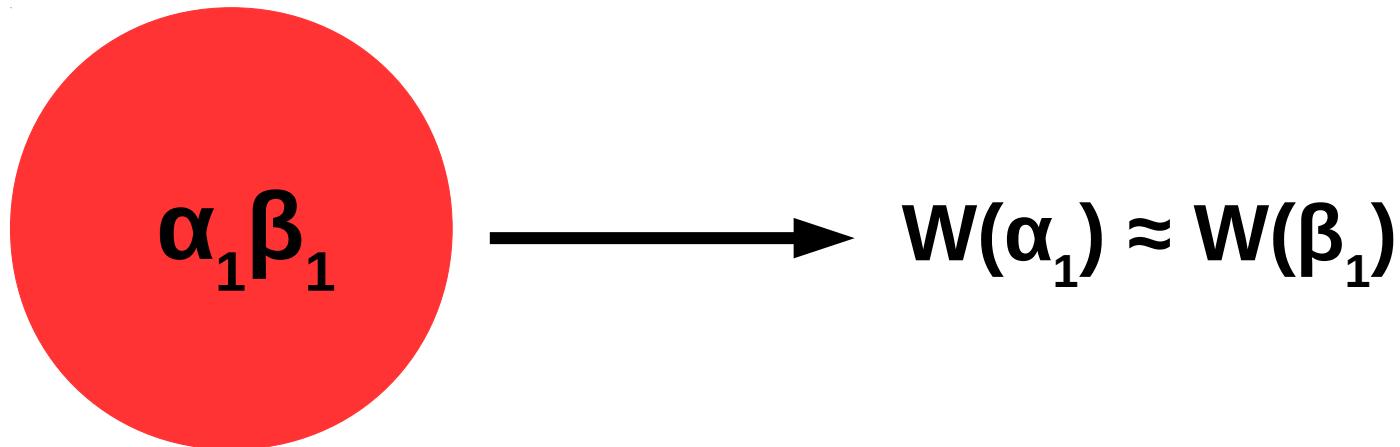
**F** Clonal frequencies are estimated with maximum-likelihood and used to distinguish  $\beta$ -sharing and dual TCR $\alpha$



**G** The output is a list of single and dual TCR clones with their respective clonal frequencies



# CHALLENGE 0: Pairing recombinants



- Expect recombinants to appear together in wells more often than they appear apart
- Pair well subsets, **not** recombinants!
  - Much fewer of the former!

# “Pair well subsets, not recombinants!”

- Minor tangent:
  - map alpha well subsets to alpha recombinants
$$A : P(\{1, 2, \dots, w\}) \rightarrow P(\{\alpha_1, \dots, \alpha_n\})$$
  - map beta well subsets to beta recombinants
$$B : P(\{1, 2, \dots, w\}) \rightarrow P(\{\beta_1, \dots, \beta_m\})$$
- If  $|A(s)|=1$ , we say the well subset  $s$  is (alpha) ***unique***, and the corresponding recombinant  $A(s)$  is ***uniquely identifiable***



# CIP Algorithm v0.0

Function **CLONEIFY**(A,B) :

    C = **PAIR**(A,B)

    Return C

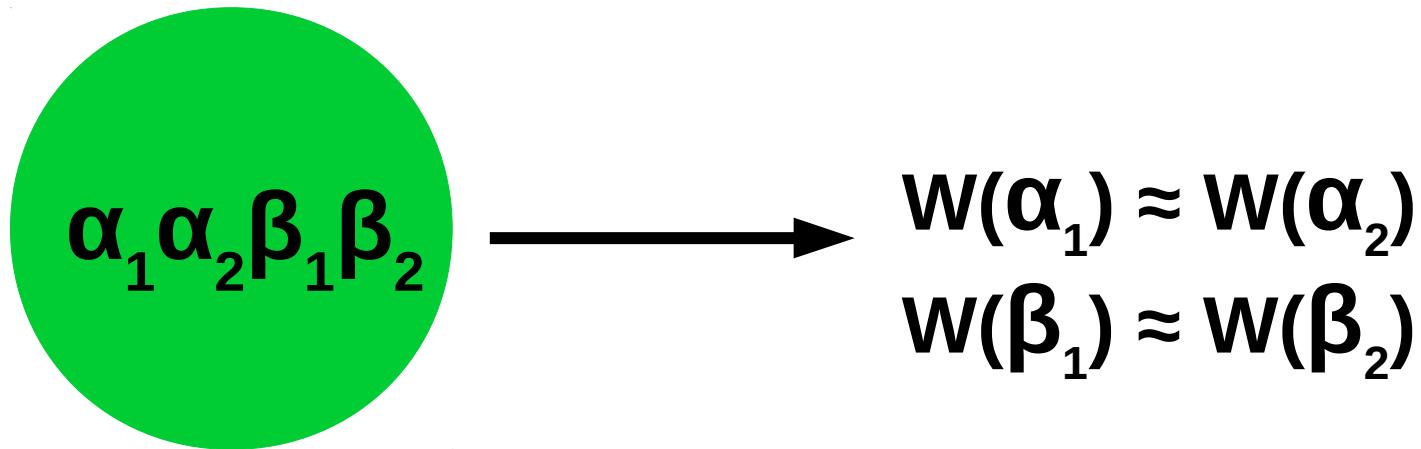
Function **PAIR**(A,B) :

    Return { (a<sub>i</sub>,b<sub>j</sub>) : a<sub>i</sub> ≈ b<sub>j</sub>, a<sub>i</sub> ∈ A, b<sub>j</sub> ∈ B }



To define later...

# CHALLENGE 1: Multiple rearrangements



- Modeling rearrangements is **easy!**
- Natural imputation of well sets

$$(a_1 \cup a_2 : (a_1, a_2) \in \text{PAIR}(A, A))$$
$$(b_1 \cup b_2 : (b_1, b_2) \in \text{PAIR}(B, B))$$


# CIP Algorithm v1.0

Function **CLONEIFY**(A, B) :

```
AD = PAIR(A, A)
BD = PAIR(B, B)
A' = AUGMENT(A, AD)
B' = AUGMENT(B, BD)
C = PAIR(A', B')
Return C
```

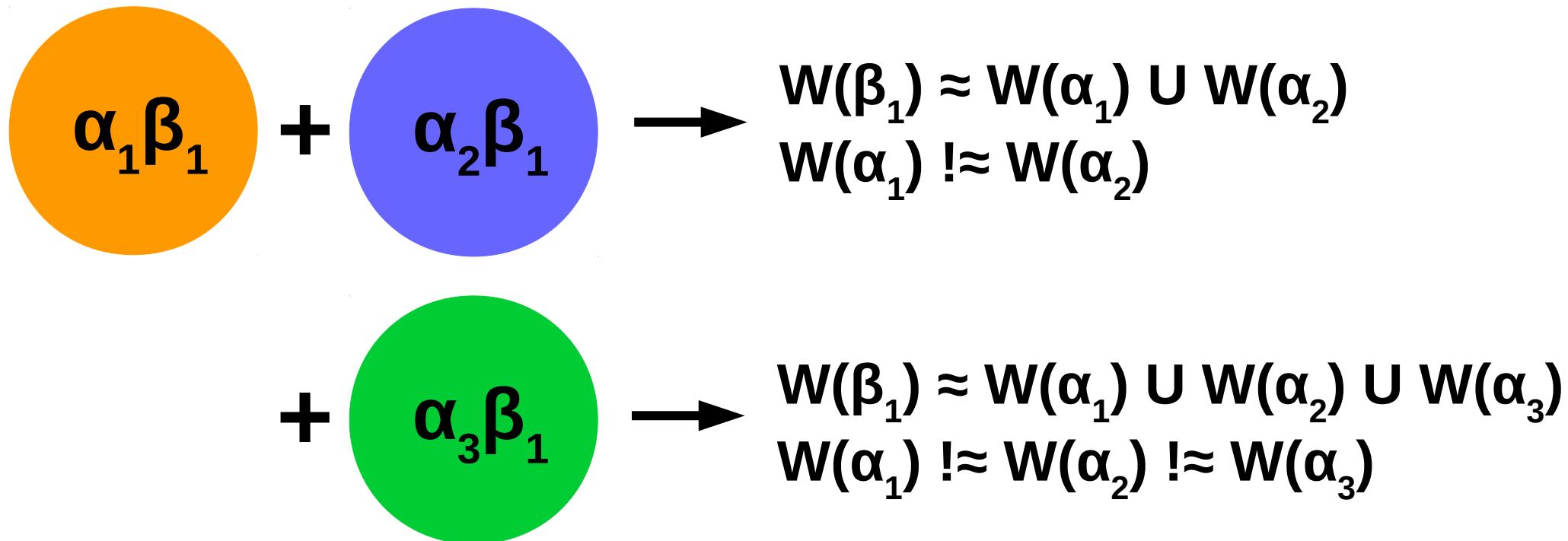
Function **AUGMENT**(W, W<sub>D</sub>) :

```
W' = W
W'((w1, w2)) = (W(w1), W(w2)), ∀(w1, w2) ∈ WD
Return W'
```

**Note:**

|W'((w<sub>1</sub>, w<sub>2</sub>))| = 1 if  
|W(w<sub>1</sub>)| = 1 and |W(w<sub>2</sub>)| = 1

# CHALLENGE 2: Shared recombinants



- Modeling sharing is *hard!*
- *Compromise:* check all pairs?
  - Cubic running time?? (not really)

# CIP Algorithm v2.0

```
Function CLONEIFY(A,B):  
    AD = PAIR(A,A)  
    BD = PAIR(B,B)  
    A' = AUGMENT(A,AD)  
    B' = AUGMENT(B,BD)  
    C = PAIR(A',B')  
    C' = SHAREIFY(A',B',C)  
    Return C'
```

```
Function SHAREIFY(A',B',C):  
    C' = C  
    For a in A':  
        For b1,b2 in C with a:  
            If b1 and b2 relatively  
            random && a ≈ b1 U b2:  
                C' = C' U {(a,b1),(a,b2)}  
    Return C'
```

# CHALLENGE 3: Computational efficiency?

- Need efficient implementation of PAIR( $A, B$ )

# CHALLENGE 3: Computational efficiency!

- Implement  $\text{PAIR}(A, B, W)$  using locality-sensitive hashing (LSH)
  - with Jaccard similarity coefficient (i.e. MinHashing)
$$SIM(a, b) = \frac{a \cap b}{a \cup b}$$
  - $a \approx b$  then means  $SIM(a, b)$  is greater than some defined threshold
  - time almost linear in the sizes of  $A$  and  $B$
- Try b-matching to put global constraint on pairing?
  - $O(e \cdot (n+e) \log(U) \log(n))$  by reducing to a min-cost max-flow and using Capacity Scaling algorithm

# *'Greedy' vs b-matching CIP Solvers*

Function **CLONEIFY**(A,B) :

```
AD = PAIR(A,A)  
BD = PAIR(B,B)  
A' = AUGMENT(A,AD)  
B' = AUGMENT(B,BD)  
C = PAIR(A',B')  
C' = SHAREIFY(A',B',C)  
Return C'
```

Function **CLONEIFY**(A,B) :

```
AD = PAIR(A,A)  
BD = PAIR(B,B)  
A' = AUGMENT(A,AD)  
B' = AUGMENT(B,BD)  
E = PAIR(A',B')  
C = BMATCH(A',B',E)  
C' = SHAREIFY(A',B',C)  
Return C'
```

**'Greedy'**

**b-matching**



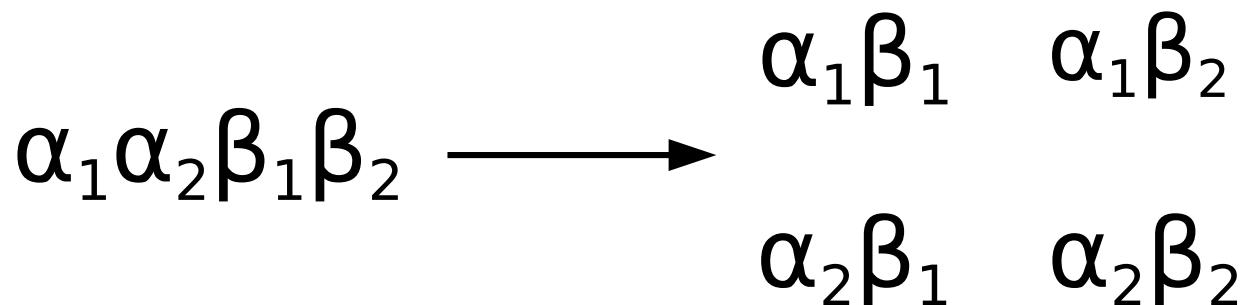
# CHALLENGE 4: Elegance?

- Kind of!
  - Simple but comprehensive model for identifying singles and duals
  - Sharing identification needs improvement...
- But how well does it all work?

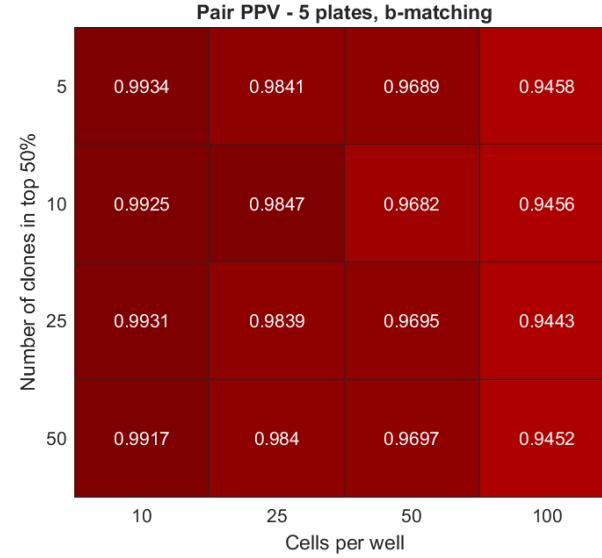
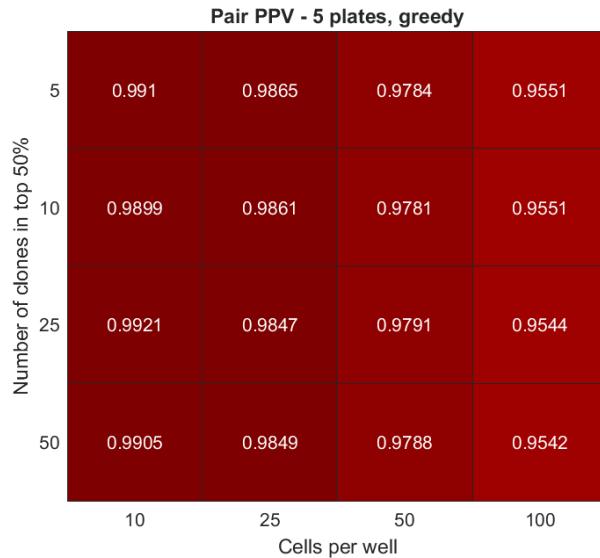
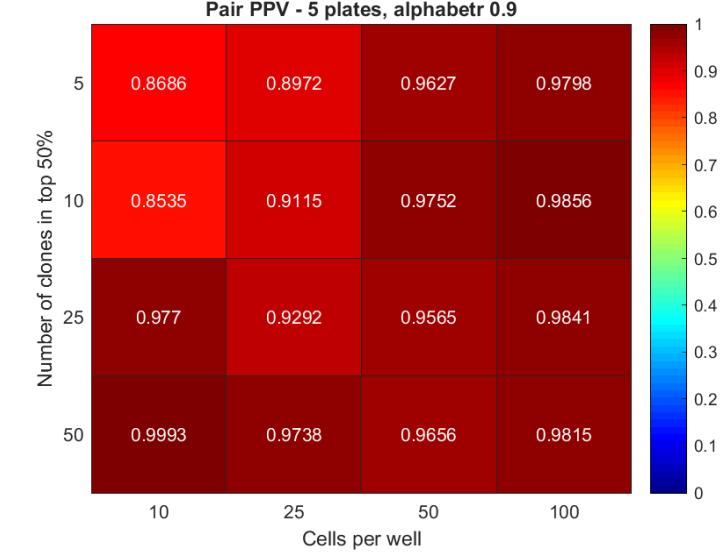
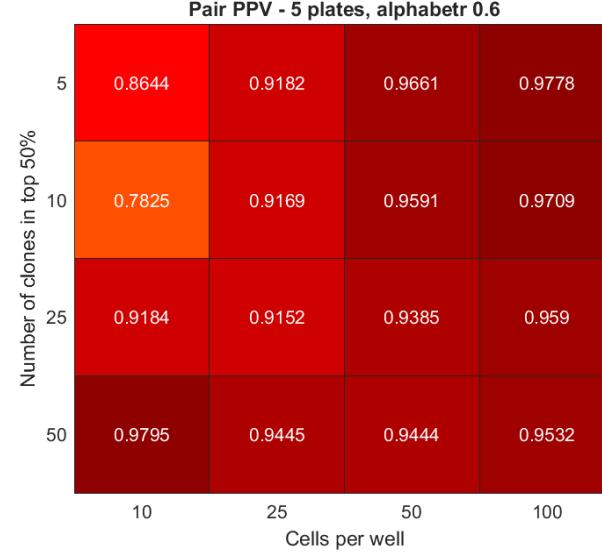
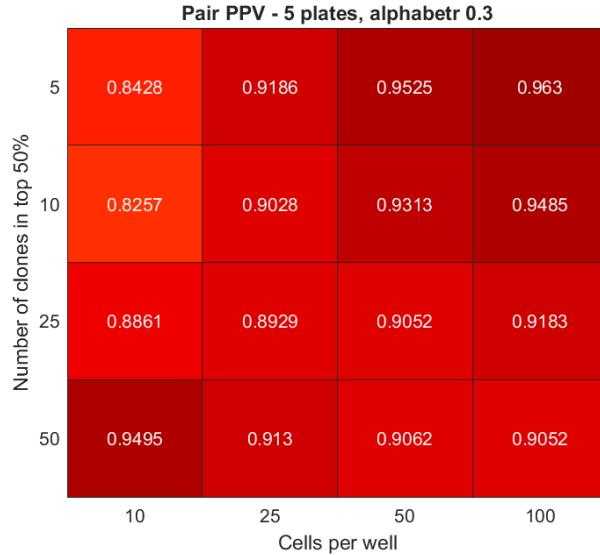


# Results

- ALPHABETR simulator
  - 2000 clones
  - 30% dual alpha, 6% dual beta
  - same sharing distribution
  - average of 50 runs
- Clone performance vs pair performance



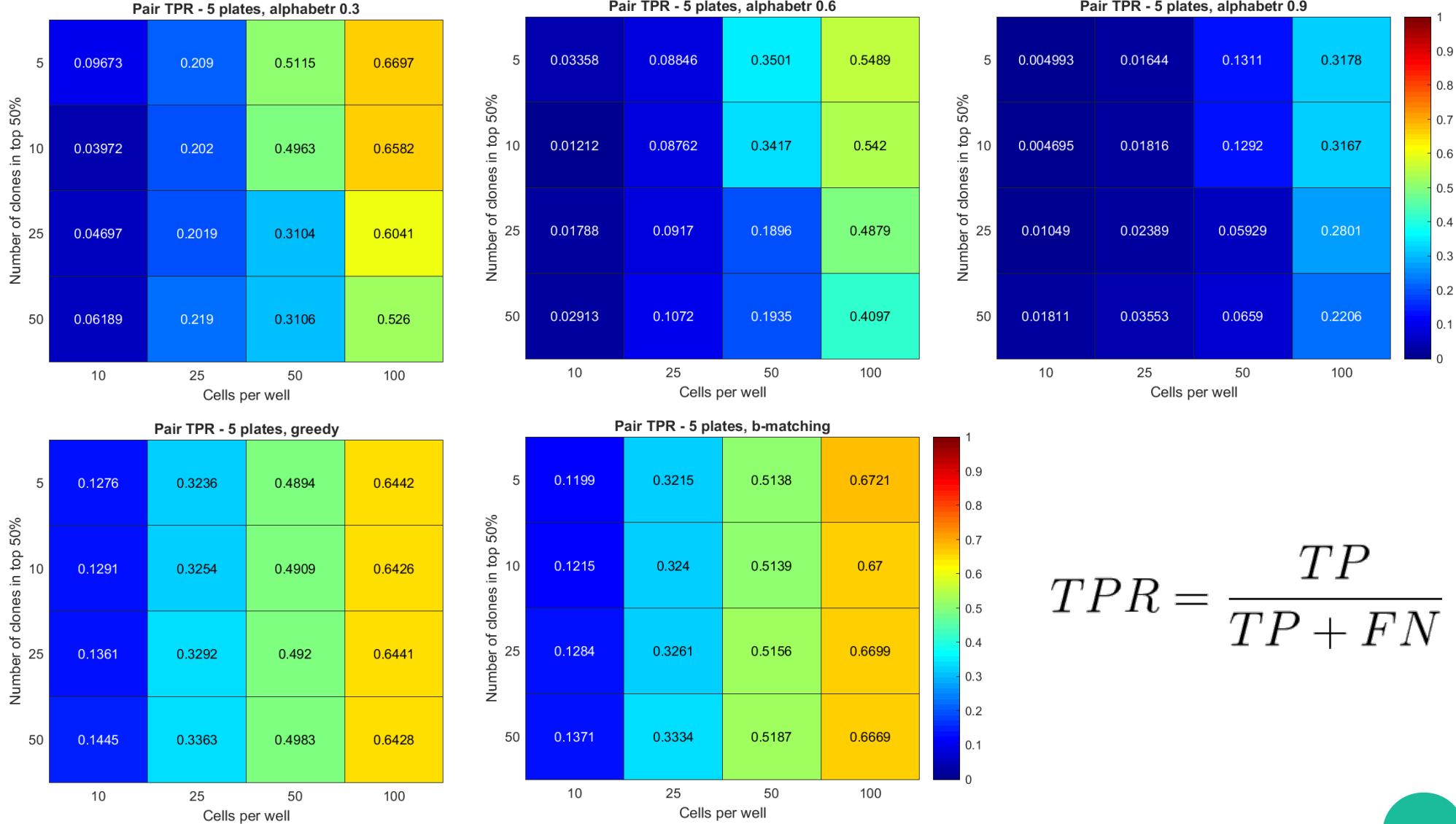
# Pair positive predictive value (PPV)



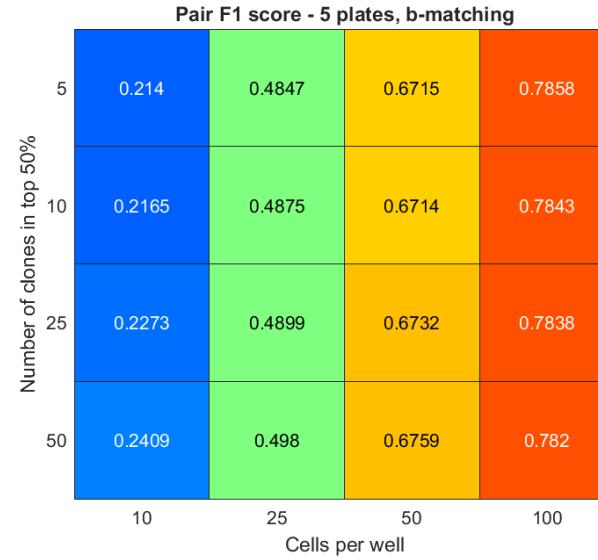
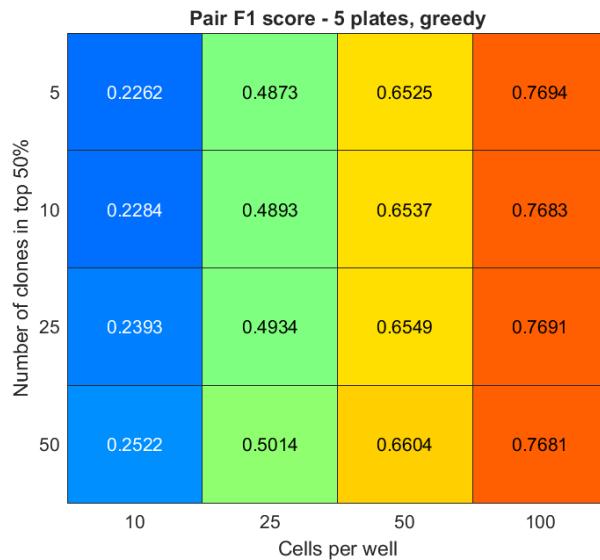
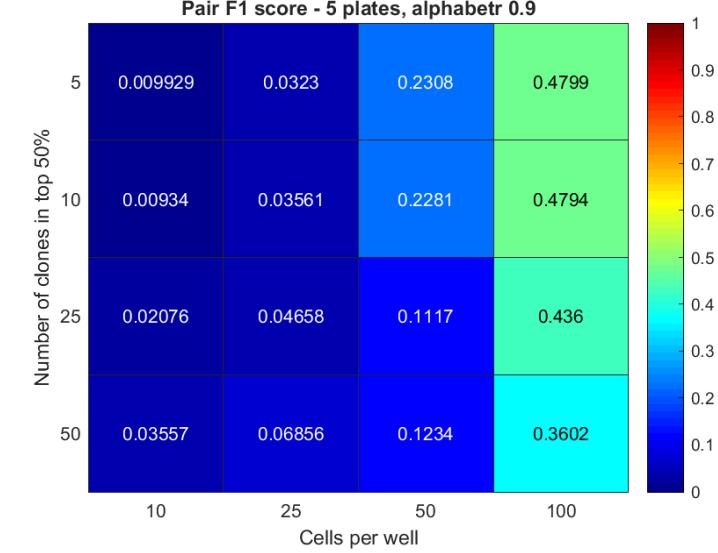
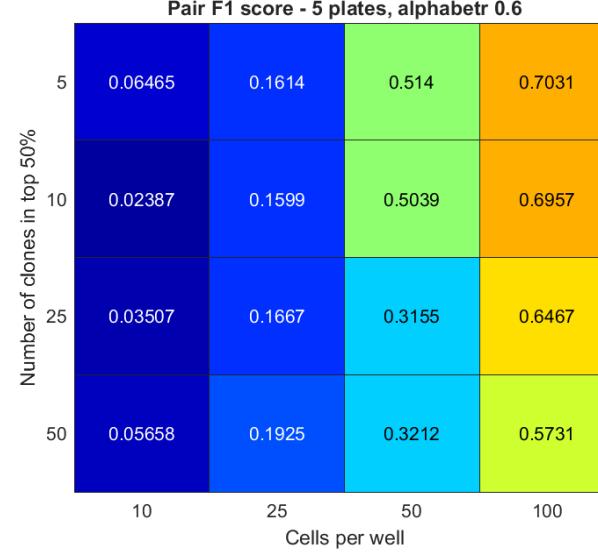
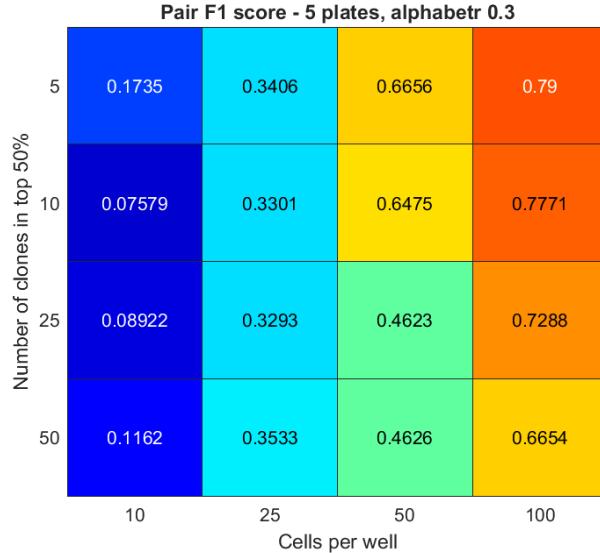
$$PPV = \frac{TP}{TP + FP}$$



# Pair true positive rate (TPR)

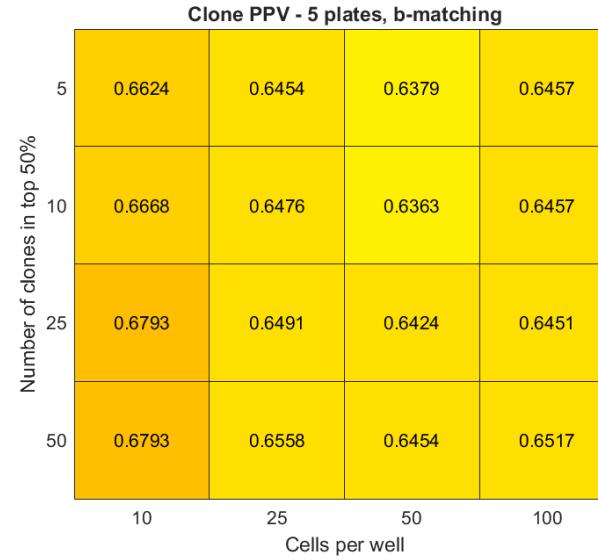
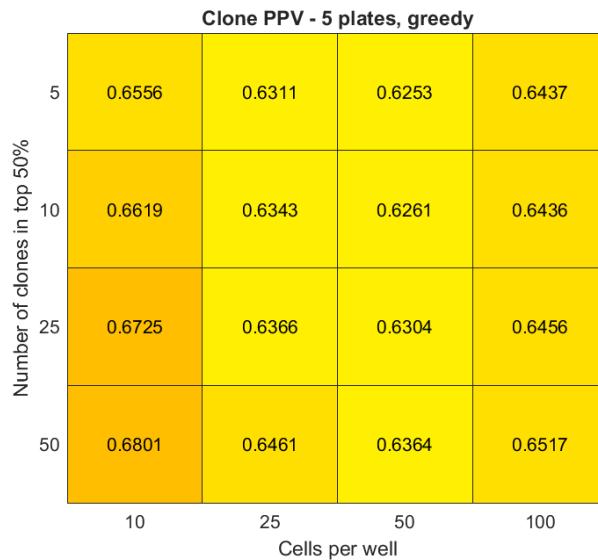
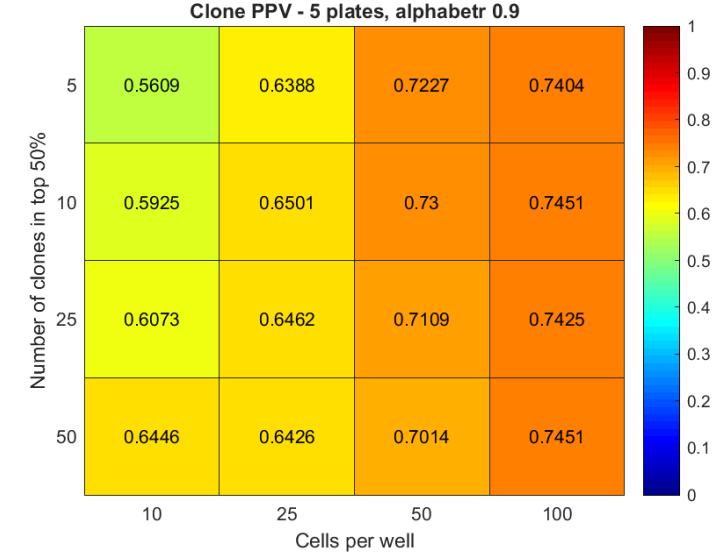
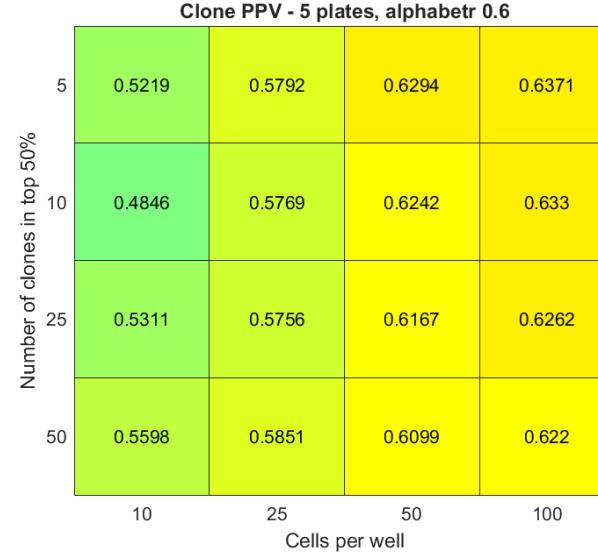
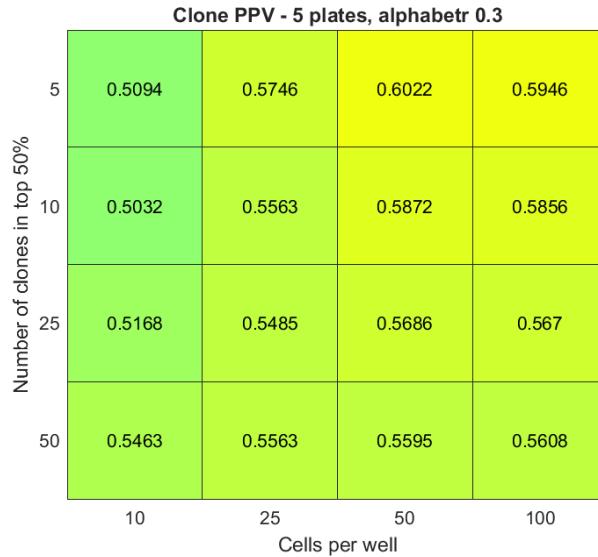


# Pair F1 score



$$F1 = \frac{2 \cdot PPV \cdot TPR}{PPV + TPR}$$

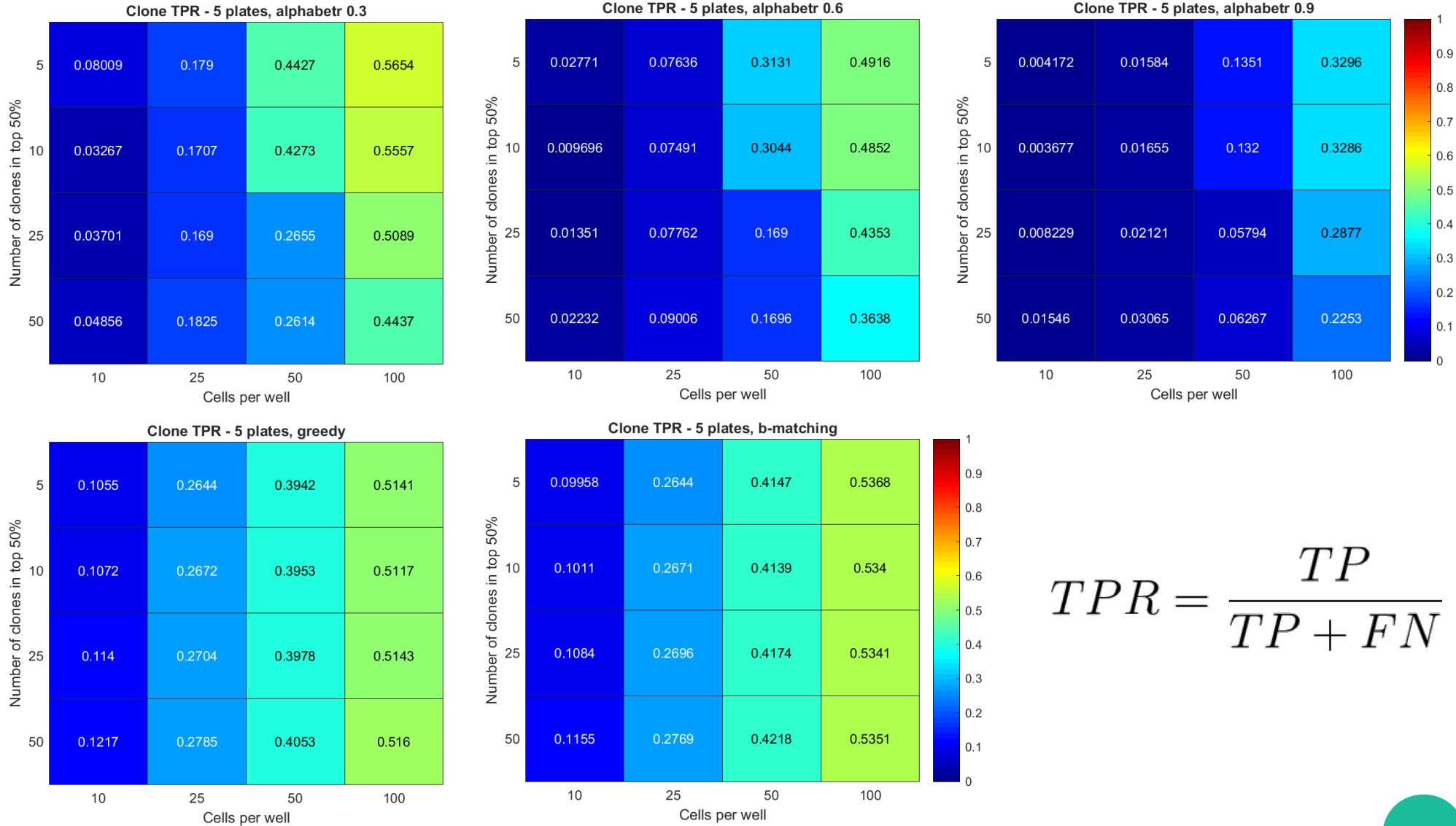
# Clone positive predictive value (PPV)



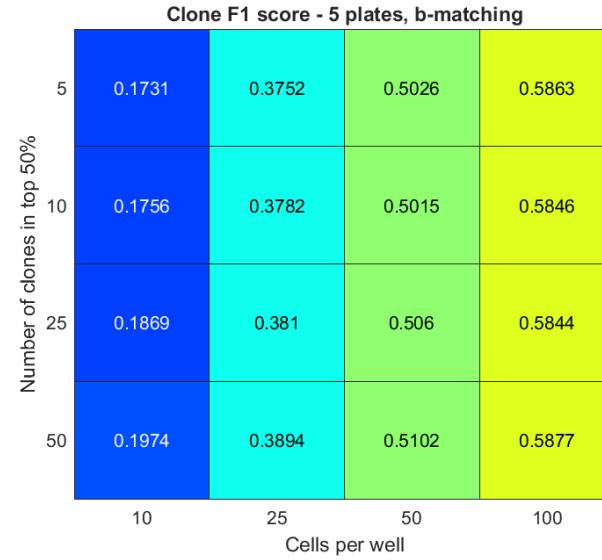
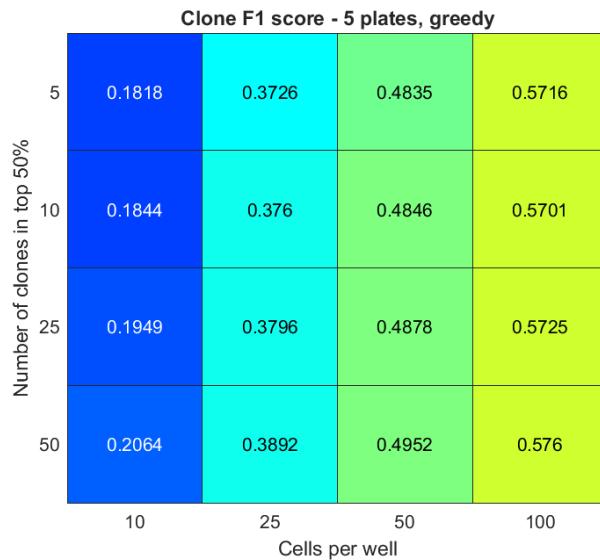
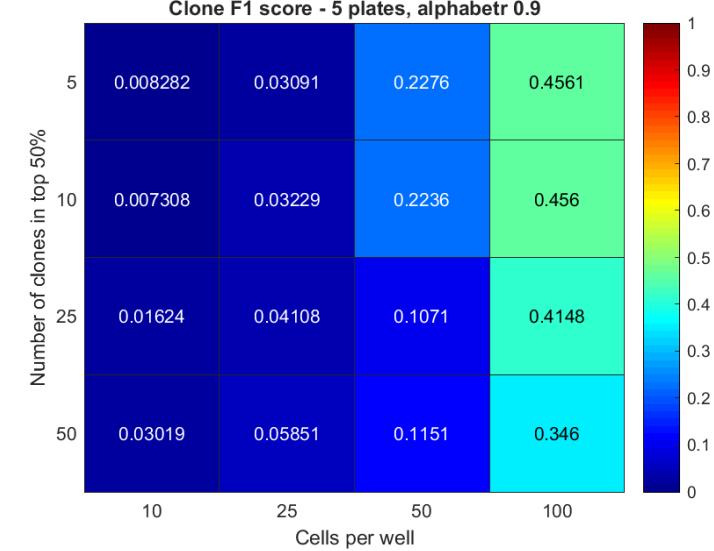
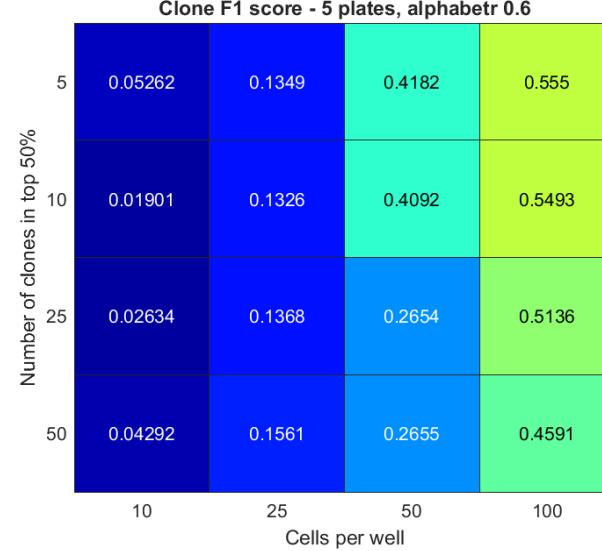
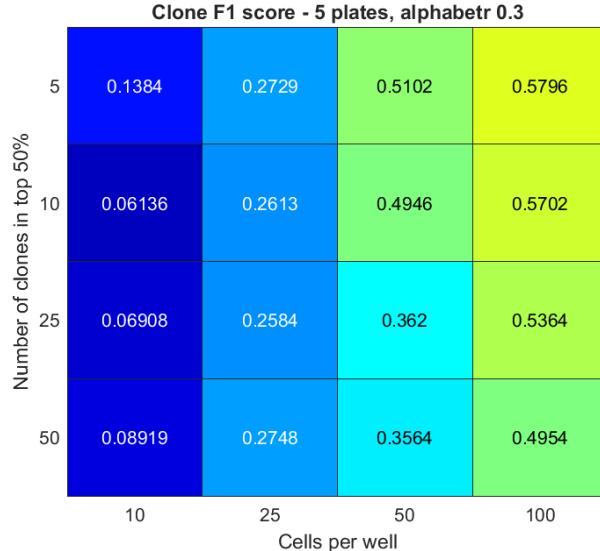
$$PPV = \frac{TP}{TP + FP}$$



# Clone true positive rate (TPR)



# Clone F1 score



$$F1 = \frac{2 \cdot PPV \cdot TPR}{PPV + TPR}$$



# Ongoing and future work

- Integrate MinHashing with solvers
- Run scaling experiments to get a concrete sense of how runtime grows
- Test solvers against real data

# Acknowledgement(s)

- Thanks Ion!



# Questions?

