

Tight Approximability Results for Test Set Problems in Bioinformatics

Piotr Berman[†]

Department of Computer Science
& Engineering
Pennsylvania State University
University Park, PA 16802
Email: berman@cse.psu.edu

Bhaskar DasGupta[‡]

Department of Computer Science
University of Illinois at Chicago
Chicago, IL 60607
Email: dasgupta@cs.uic.edu

Ming-Yang Kao[¶]

Department of Computer Science
Northwestern University
Evanston, IL 60201
Email: kao@cs.northwestern.edu

[†] Supported by NSF grant CCR-0208821.

[‡] Supported by NSF Grants CCR-0296041, CCR-0206795, CCR-0208749 and a CAREER grant IIS-0346973.

[¶] Supported by NSF grant EIA-0112934.

A General Framework

Problem $\text{TS}^\Gamma(k)$: $\Gamma \subseteq 2^{\{0,1,2\}}$, k a positive integer k

Instance: (n, \mathcal{S}) where $\mathcal{S} \subseteq 2^{\{0,1,2,\dots,n-1\}}$

Terminologies:

- A **k-test** is a union of at most k sets from \mathcal{S}
- For a $\gamma \in \Gamma$ and two distinct elements $x, y \in \{0, 1, 2, \dots, n-1\}$, a **k-test T γ -distinguishes x and y** if $|\{x, y\} \cap T| \in \gamma$.

Valid solutions: A collection \mathcal{T} of k -tests such that

$$(\forall x, y \in \{0, 1, 2, \dots, n-1\} \quad \forall \gamma \in \Gamma) \quad x \neq y$$

$$\implies$$

$$\exists T \in \mathcal{T} \text{ such that } T \text{ } \gamma\text{-distinguishes } x \text{ and } y$$

Objective: minimize $|\mathcal{T}|$.

Some Problems Captured by General Framework

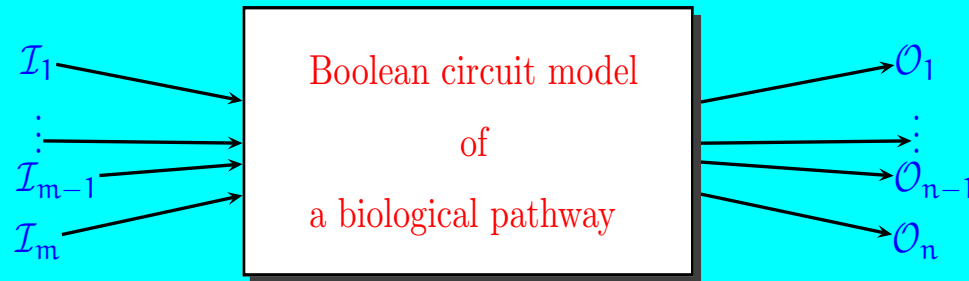
Minimum Test Collection Problems

- Equivalent to $TS^{\{1\}}(1)$
 - Objects $\mathcal{O}_1, \mathcal{O}_2, \dots, \mathcal{O}_n$
 - $\mathcal{O}_i, \mathcal{O}_j$ is **distinguished** by test T if $T \cap \{\mathcal{O}_i, \mathcal{O}_j\} = 1$.
- Applications: diagnostic testing
- Reference: Garey and Johnson's book on NP-completeness, page 71.

More Problems Captured by General Framework

Condition Cover Problems

- Captured by $\text{TS}^{\{1\},\{0,2\}}(1)$

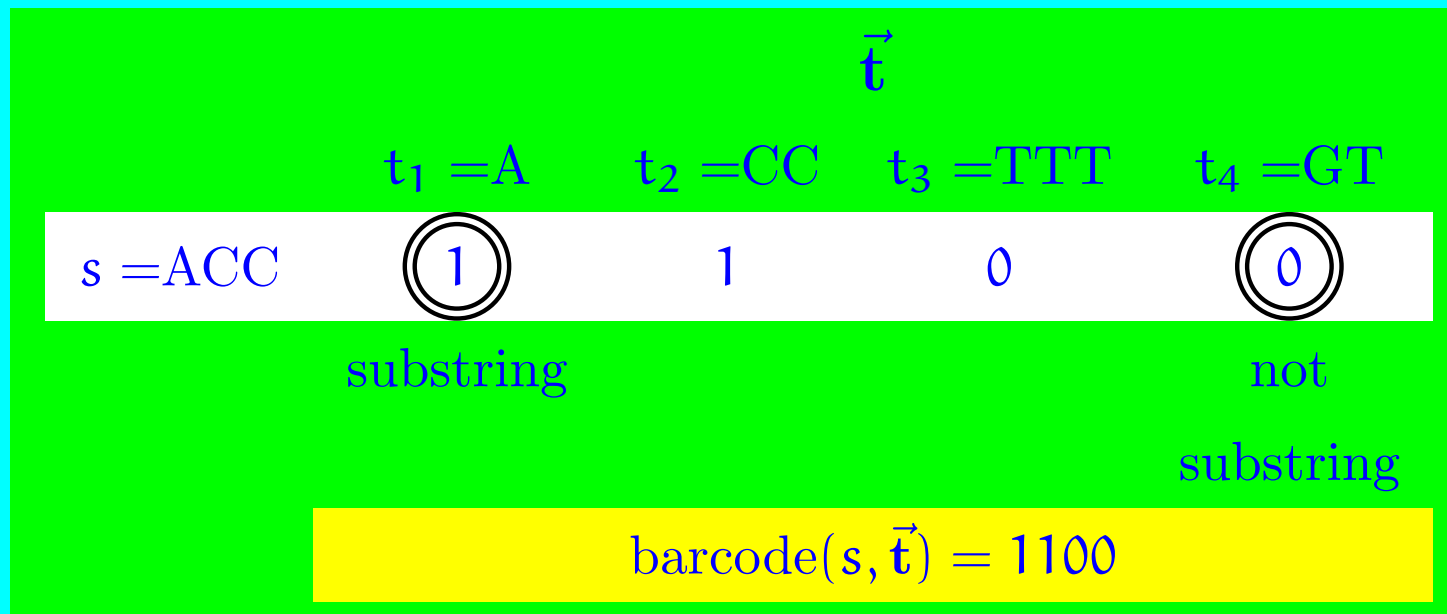


- **test** : assignment σ of I_1, I_2, \dots, I_m to 0/1 values.
For this assignment, m^{th} output is $O_m(\sigma)$
- Set of tests given as part of input
- O_i, O_j are **distinguished** if
 - $\exists \sigma : O_i(\sigma) = O_j(\sigma) \ \& \ \exists \rho : O_i(\rho) \neq O_j(\rho)$
- **Applications**: verifying a multi-output feedforward Boolean circuit as a model of specific **biological pathways**

More Problems Captured by General Framework

Simplest String Barcoding Problems ($SB^\Sigma(1)$)

- very special case of $TS^{\{1\}}(k)$
 - **Given:** set \mathcal{S} of sequences over alphabet Σ
 - **Definition of barcode:** for a sequence s and a set of sequences \mathbf{t} , $\text{barcode}(s, \vec{\mathbf{t}})$ is the Boolean vector (c_0, c_1, c_{m-1}) where c_i is 1 if t_i is a **substring** of s .



Simplest String Barcoding Problems (continued)

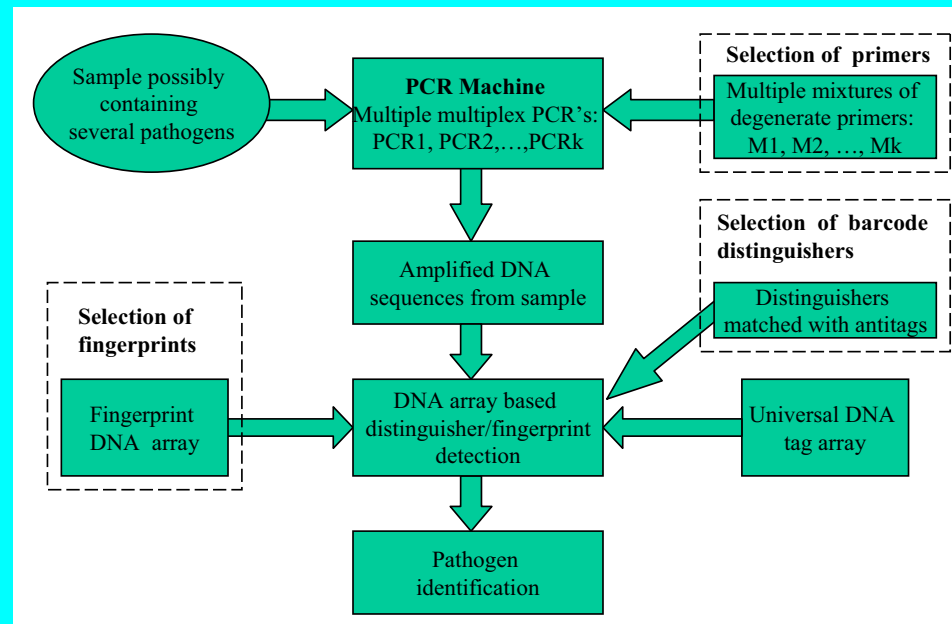
- **Valid solutions:** a set of sequences \vec{t} such that
 $\forall s, s' \in \mathcal{S} : s \neq s' \equiv \text{barcode}(s, \vec{t}) \neq \text{barcode}(s', \vec{t})$

\mathcal{S}	\vec{t}			
	A	CC	TTT	GT
$S_1 = \text{AAC}$	1	0	0	0
$S_2 = \text{ACC}$	1	1	0	0
$S_3 = \text{GGGG}$	0	0	0	0
$S_4 = \text{GTGTGG}$	0	0	0	1
$S_5 = \text{TTTT}$	0	0	1	0

Objective: minimize $|\vec{t}|$.

Simplest String Barcoding Problems (continued)

- Applications:
 - database compression/fast database search for DNA sequences
 - DNA microarray designs for unknown pathogen identification



More Problems Captured by General Framework

Minimum Cost Probe Set with Threshold r ($\text{MCP}^\Sigma(r)$):

- variation of $\text{TS}^{\{1\}}(1)$

Given : sets \mathcal{S} and \mathcal{P} of sequences over alphabet Σ and an integer $r > 0$

Definition of r -barcode : for a sequence s and a set of sequences \mathbf{t} , r -barcode($s, \vec{\mathbf{t}}$) is the integer vector (c_0, c_1, c_{m-1}) where

$$c_i = \min \{ r, \text{number of occurrences of } t_i \text{ in } s \}$$

Example of 2-barcode ($r = 2$)

\vec{t}

$t_1 = A$ $t_2 = CC$ $t_3 = AC$ $t_4 = G$

$s = ACCCCA$ $\textcircled{2}$ $\textcircled{2}$ 1 0

$\min\{2, 2\}$ $\min\{2, 3\}$
 r r

2-barcode(s, \vec{t}) = 2210

Valid solutions: set of sequences $\vec{t} \subseteq \mathcal{P}$ such that

$$\forall s, s' \in \mathcal{S} : s \neq s' \Rightarrow r\text{-barcode}(s, \vec{t}) \neq r\text{-barcode}(s', \vec{t})$$

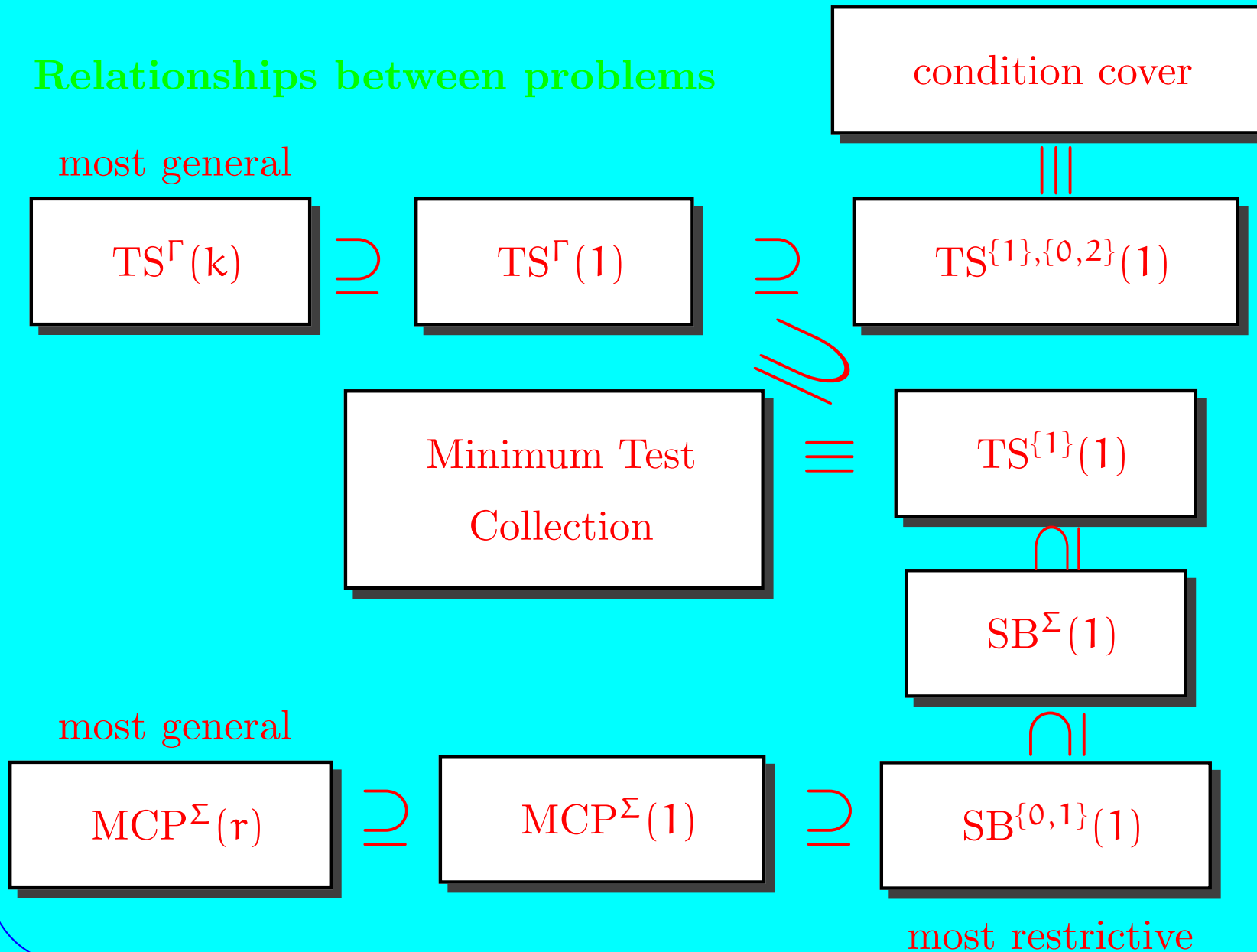
Objective: minimize $|\vec{t}|$.

- A special case of $\text{MCP}^{\Sigma}(r)$ is $\text{SB}^{\Sigma}(1)$

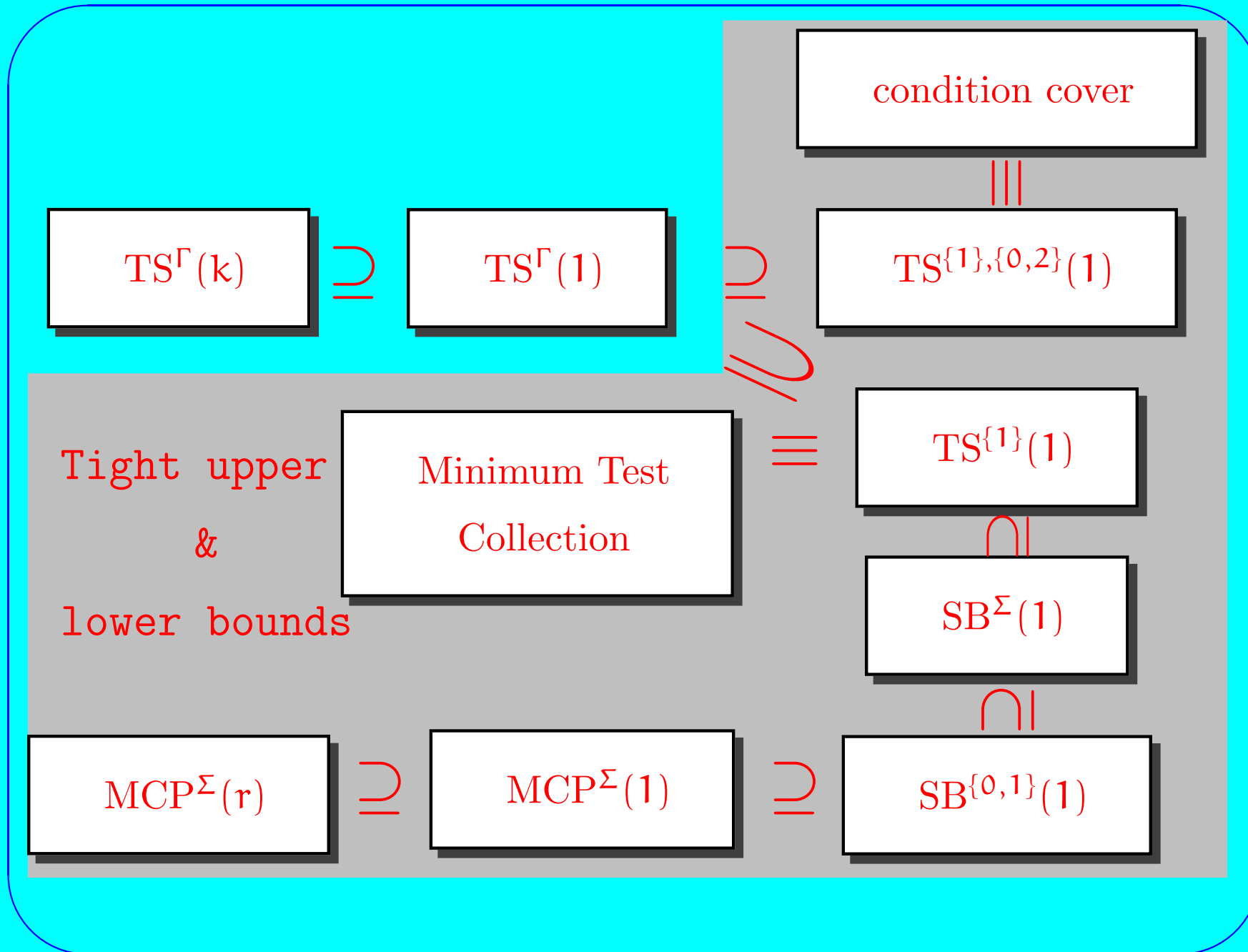
If \mathcal{P} is the set of all substrings of all sequences in \mathcal{S} then $\text{MCP}^{\Sigma}(1)$ is precisely $\text{SB}^{\Sigma}(1)$

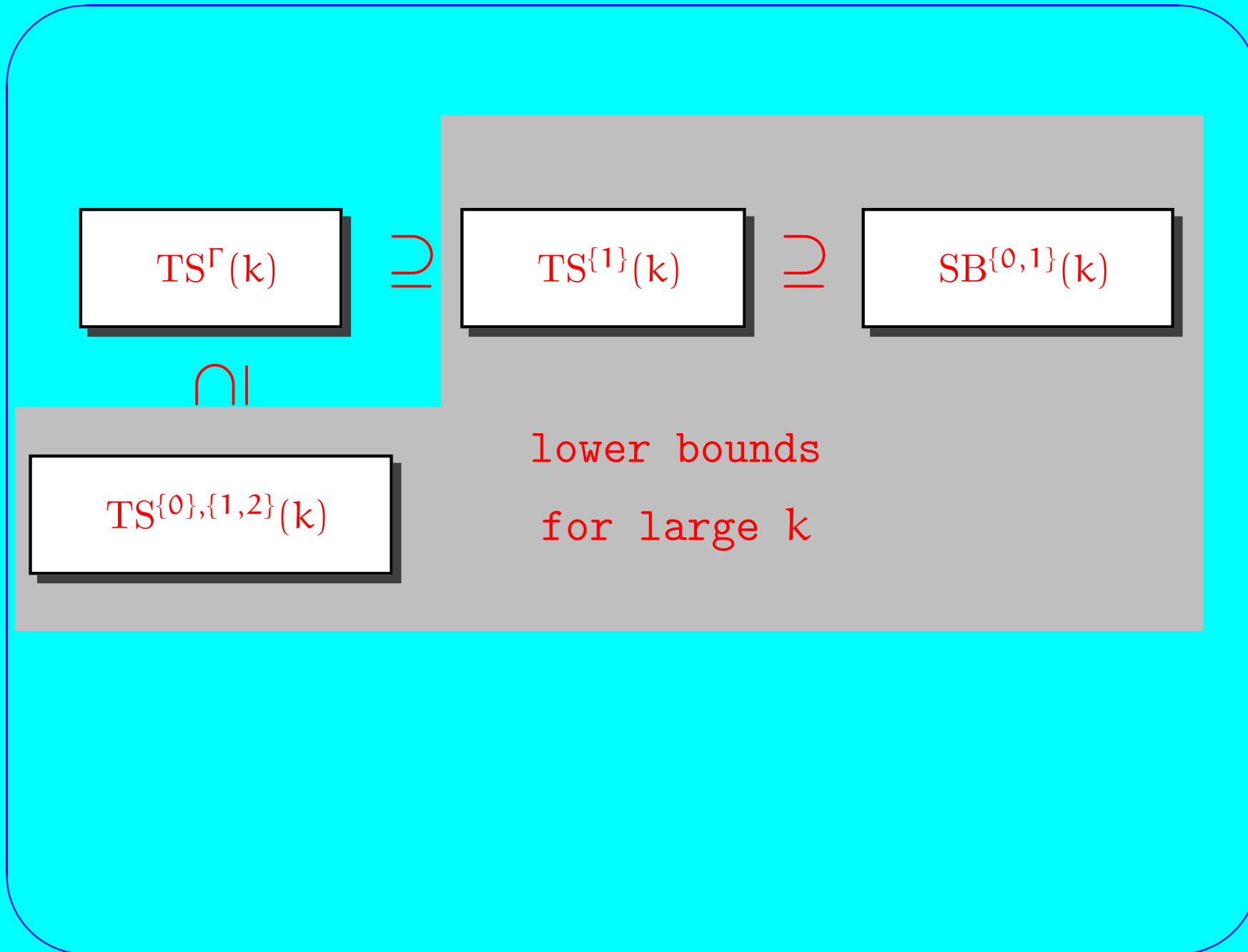
- **Applications:** in minimization the number of **oligonucleotide probes** needed for analyzing populations of ribosomal RNA gene (rDNA) clones by **hybridization experiments on DNA microarrays**

Relationships between problems



- Upper bounds are proved for the more general cases ($\text{TS}^{\{1\},\{0,2\}}(1)$, $\text{TS}^{\{1\}}(1)$ and $\text{MCP}^{\Sigma}(r)$)
- Lower bounds are proved for the most restrictive case ($\text{SB}^{\{0,1\}}(1)$ and $\text{TS}^{\{1\},\{0,2\}}(1)$)
- Lower bounds are also proved for $\text{TS}^{\{1\},\{0,2\}}(k)$, $\text{TS}^{\{1\}}(k)$ and $\text{SB}^{\{0,1\}}(k)$ when k is large





Summary of our results

(matching upper/lower bounds)

Problem	Approximation Ratio		
	Upper Bound (algorithm)	Lower Bound	
		the bound	Assumptions
$TS^{\{1\}}(1)$	$1 + \ln n$	$(1 - \varepsilon) \ln n$	$NP \not\subseteq DTIME(n^{\log \log n})$
$TS^{\{1\},\{0,2\}}(1)$	$1 + \ln 2 + \ln n$	$(1 - \varepsilon) \ln n$	$NP \not\subseteq DTIME(n^{\log \log n})$
$SB^{\Sigma}(1)$	$1 + \ln n$	$(1 - \varepsilon) \ln n$	$NP \not\subseteq DTIME(n^{\log \log n})$ $ \Sigma > 1$
$MCP^{\Sigma}(r)$	$[1 + o(1)] \ln n$	$(1 - \varepsilon) \ln n$	$NP \not\subseteq DTIME(n^{\log \log n})$ $ \Sigma > 1$
$TS^{\{1\}}(n^{\delta})$ $TS^{\{1\},\{0,2\}}(n^{\delta})$		n^{ε}	$NP \neq co-RP$ $0 < \varepsilon < \delta < 1$
$SB^{\{0,1\}}(n^{\delta})$		n^{ε}	$NP \neq co-RP$ $0 < \varepsilon < \delta < \frac{1}{2}$

ε and δ are arbitrary constants

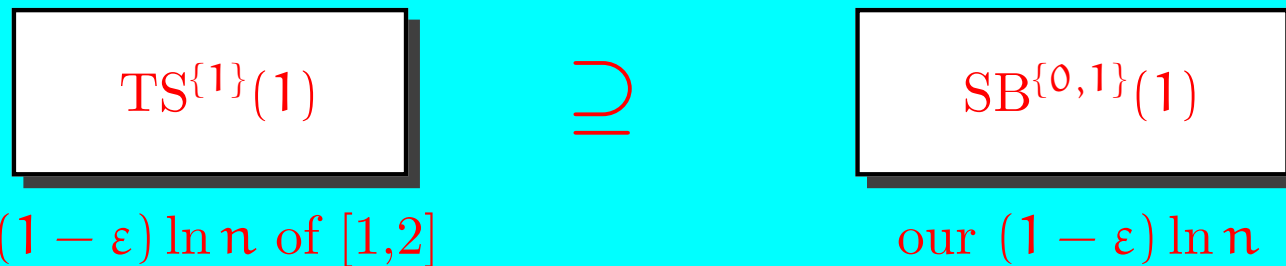
Comparison of our results with those in [1,2]

- [1] B. V. Halldórsson, M. M. Halldórsson and R. Ravi. *On the Approximability of the Minimum Test Collection Problem*, Proc. Ninth Annual European Symposium on Algorithms, Lecture Notes in Computer Science 2161, pp. 158-169, 2001.
- [2] K. M. J. De Bontridder, B. V. Halldórsson, M. M. Halldórsson, C. A. J. Hurkens, J. K. Lenstra, R. Ravi and L. Stougie. *Approximation algorithms for the test cover problem*, Mathematical Programming-B, Vol. 98, No. 1-3, 2003, pp. 477-491.

Comparison of our results with those in [1,2]

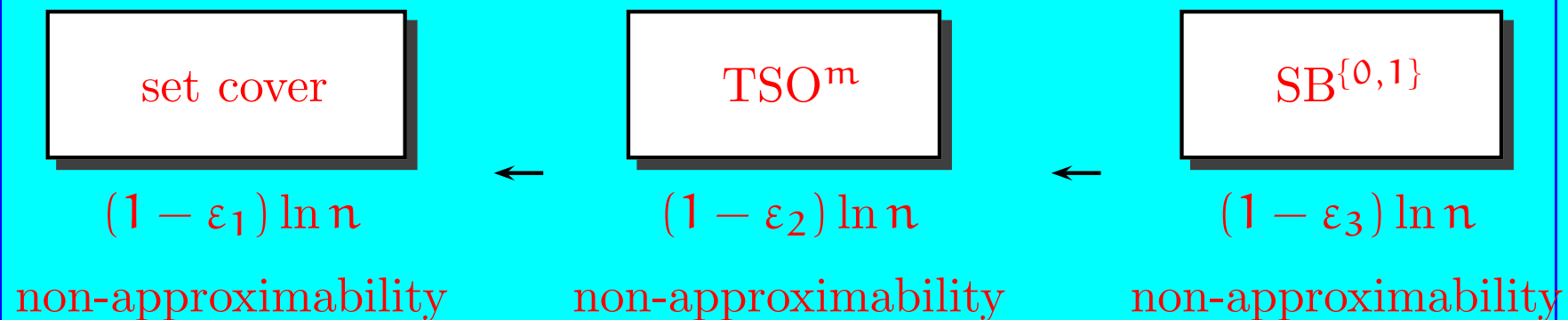
- The authors in [1,2] proved a $(1 - \varepsilon) \ln n$ lower bound for approximation for $TS^{\{1\}}(1)$

We prove a lower bound of $(1 - \varepsilon) \ln n$ for the very special case $SB^{\{0,1\}}$



Comparison of our results with those in [1,2] (continued)

- The proof in [1,2] from set-cover to $TS^1(1)$ does not seem to be easily transformable to provide a lower bound for $SB^{\{0,1\}}$ with a similar quality of non-approximability because of the special nature of $SB^{\{0,1\}}$
- We therefore needed to introduce an artificial intermediate problem (the “test set with order with positive integer parameter m ” problem, denoted by TSO^m) which we could then translate to $SB^{\{0,1\}}$ in a non-trivial manner



Comparison of our results with those in [1,2] (continued)

- In general, TSO^m is neither equivalent to or nor a special case of $\text{TS}^{\{1\}}(1)$.

Summary of Other Techniques Used

- Algorithm for $TS^{\{1\}}(1)$, $TS^{\{1\},\{0,2\}}(1)$ and $MCP^{\Sigma}(r)$ is a greedy algorithm that selects tests based on the change of **information content** of the partition of the universe

- A set of tests \mathcal{T} defines an **entropy** $H_{\mathcal{T}}$
- notion of **information content** IC

\mathcal{T} = already selected tests

T = new test

$$IC(T, \mathcal{T}) = H_{\mathcal{T}} - H_{\mathcal{T} \cup T}$$

- Greedy heuristics

$\mathcal{T} = \emptyset$

while $H_{\mathcal{T}} \neq 0$ **do**

 select a $T \in \mathcal{S} - \mathcal{T}$ that *maximizes* $IC(T, \mathcal{T})$

$\mathcal{T} = \mathcal{T} \cup T$

endwhile

- The inapproximability results for $\text{TS}^{\{1\}}(\mathfrak{n}^\delta)$, $\text{TS}^{\{0\},\{1,2\}}(\mathfrak{n}^\delta)$ and $\text{SB}^{\{0,1\}}(\mathfrak{n}^\delta)$ are obtained by **approximation preserving reductions** from the **graph coloring** problem.

Thank You for your attention!!