

Optimization in the Race to a Liquid Biopsy

Kyra Gan
Cornell Tech

Su Jia
Cornell

Andrew Li
CMU

Sridhar Tayur
CMU

A Little Biology

- Cancer is caused by DNA mutations
- Tumors contain mutated DNA

A Little Biology

- Cancer is caused by DNA mutations
 - Tumors contain mutated DNA
- **Cell-free DNA:** blood contains tiny amounts of mutated DNA
 - Concentration of one in ten-thousand
 - Should hypothetically be a cancer signal

The Race is On

- Recent academic successes based on this idea:

[Cohen et al.] Detection and localization of surgically resectable cancers with a multi-analyte blood test. *Science* (2018).

[Liu et al.] Sensitive and specific multi-cancer detection and localization using methylation signatures in cell-free DNA. *Annals of Oncology* (2020).

The Race is On

- Recent academic successes based on this idea:

[Cohen et al.] Detection and localization of surgically resectable cancers with a multi-analyte blood test. *Science* (2018).

Forbes

Jan 18, 2018, 02:00pm EST

A New \$500 Blood Test Could Detect Cancer Before Symptoms Develop

The Race is On

- Recent academic successes based on this idea:

[Cohen et al.] Detection and localization of surgically resectable cancers with a multi-analyte blood test. *Science* (2018).



- ...and biotech firms building on it:

GRAIL

 GUARDANT™

EXACT
SCIENCES

freonome

Nearing the Finish Line

- Grail's *Galleri* test

Nearing the Finish Line

- Grail's *Galleri* test
 - \$949

The cost of the Galleri[®] test may vary depending on the healthcare practice or provider who orders the test. The list price for the Galleri test is \$949.

Nearing the Finish Line

- Grail's *Galleri* test
 - \$949
 - Detects 50+ cancer types
 - **Specificity** 99.5%

50+
cancer types

99.5%
specificity¹

Nearing the Finish Line

- Grail's *Galleri* test
 - \$949
 - Detects 50+ cancer types
 - **Specificity** 99.5%
 - **Sensitivity** 77%

Sensitivity

76.3% sensitivity in cancers that cause two-thirds of cancer deaths in the US ⓘ ^{1,6,7}

Data Available Today

- DNA is 3 billion **addresses** long:

... AGCATGCAGTACGTACGTCACATTCGATCGATGG...

Data Available Today

- DNA is 3 billion **addresses** long:

... AGCATGCAGTACGTACGTACACATTTCGATCGATGG...

- **Mutations** are with respect to a reference sequence

... AG**G**ATGCAGT**C**CGTACGTACACATT**C**AATCGATGG...

Data Available Today

- DNA is 3 billion **addresses** long:

... AGCATGCAGTACGTACGTACACATTTCGATCGATGG...

- **Mutations** are with respect to a reference sequence

...00**1**00000000**1**0000000000000000**1**00000000...

Data Available Today

- DNA is 3 billion **addresses** long:

... AGCATGCAGTACGTACGTACACATTTCGATCGATGG...

- **Mutations** are with respect to a reference sequence

...00**1**00000000**1**0000000000000000**1**00000000...

- Many (hundreds of thousands) of tumors have been **sequenced**

Data Available Today

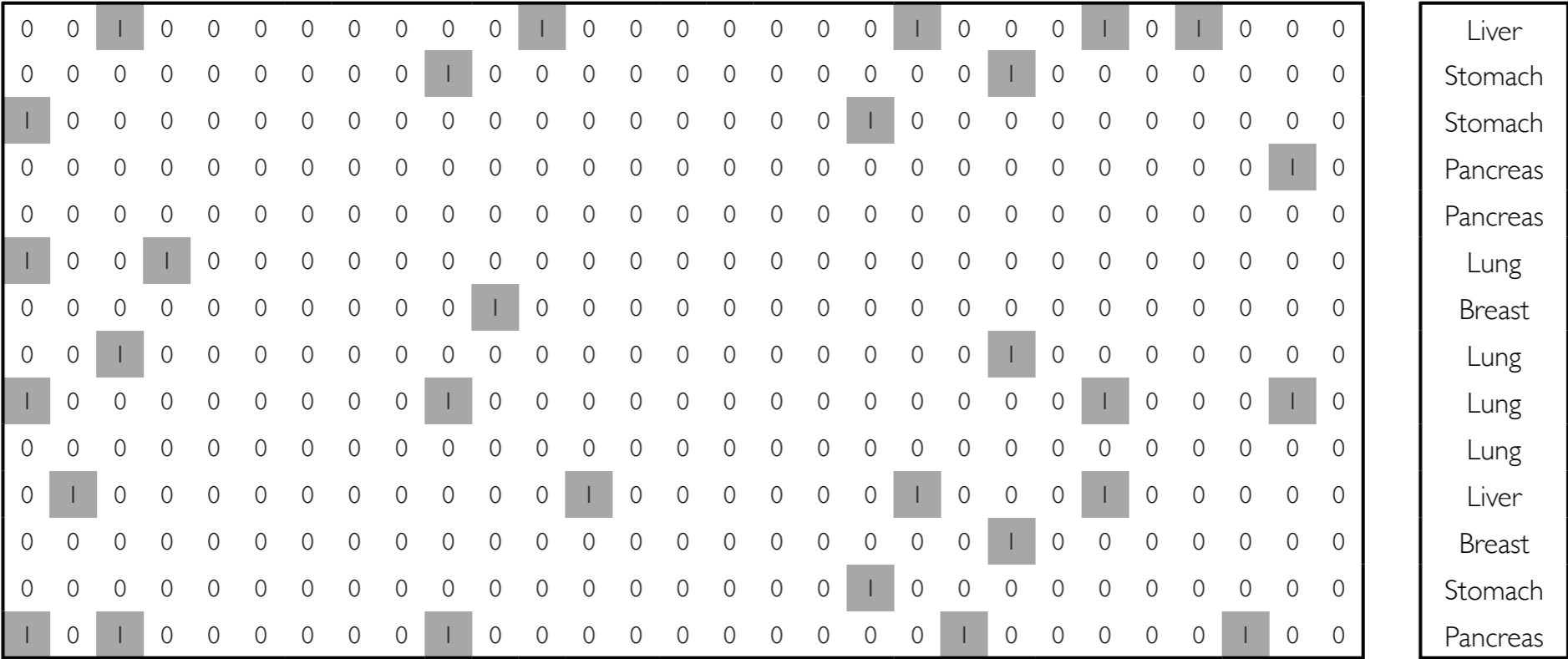
- DNA is 3 billion **addresses** long:

... AGCATGCAGTACGTACGTACACATTTCGATCGATGG...

- **Mutations** are with respect to a reference sequence

...00**1**0000000**1**0000000000000000**1**00000000...

- Many (hundreds of thousands) of tumors have been **sequenced**



Challenge: Cost

- A back-of-the-envelope calculation:
- Test must cost at most $\$10^2$

0	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	1	0	0	0	Liver
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	Stomach
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	Stomach
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	Pancreas
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Pancreas
1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Lung
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Breast
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	Lung
1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	Lung
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Lung
0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	Liver
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	Breast
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	Stomach
1	0	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	Pancreas

Challenge: Cost

- A back-of-the-envelope calculation:
 - Test must cost at most $\$10^2$
 - Sequencing an address **ten-thousand times** costs $\$10^{-2}$

0	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	1	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
1	0	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0

Liver
Stomach
Stomach
Pancreas
Pancreas
Lung
Breast
Lung
Lung
Lung
Liver
Breast
Stomach
Pancreas

Challenge: Cost

- A back-of-the-envelope calculation:
 - Test must cost at most $\$10^2$
 - Sequencing an address **ten-thousand times** costs $\$10^{-2}$
 - Thus, can only use a **panel** of $\sim 10^4$ addresses

0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	1	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
1	0	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0

Liver
Stomach
Stomach
Pancreas
Pancreas
Lung
Breast
Lung
Lung
Lung
Liver
Breast
Stomach
Pancreas

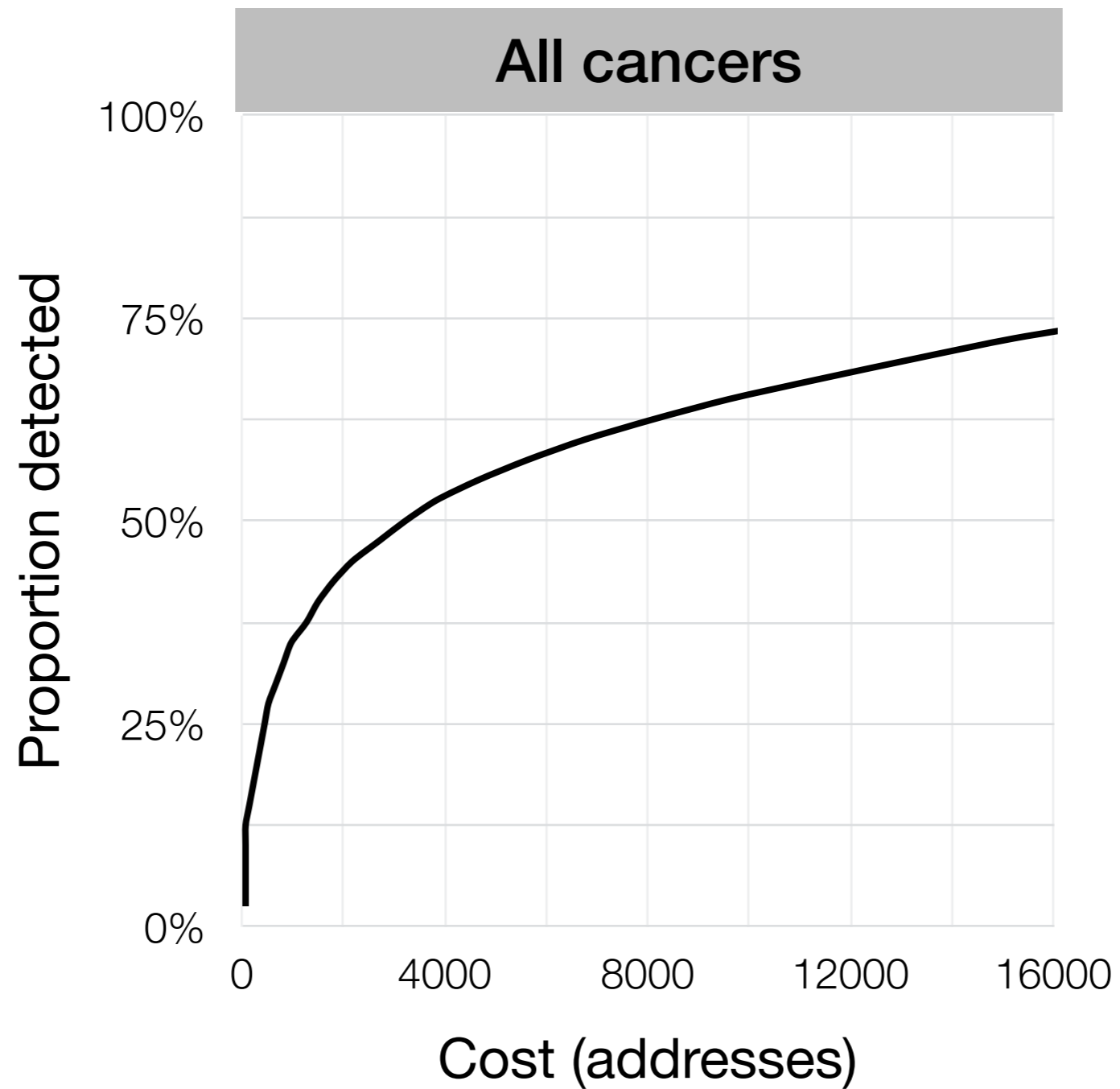
Challenge: Cost

- A back-of-the-envelope calculation:
 - Test must cost at most $\$10^2$
 - Sequencing an address **ten-thousand times** costs $\$10^{-2}$
 - Thus, can only use a **panel** of $\sim 10^4$ addresses

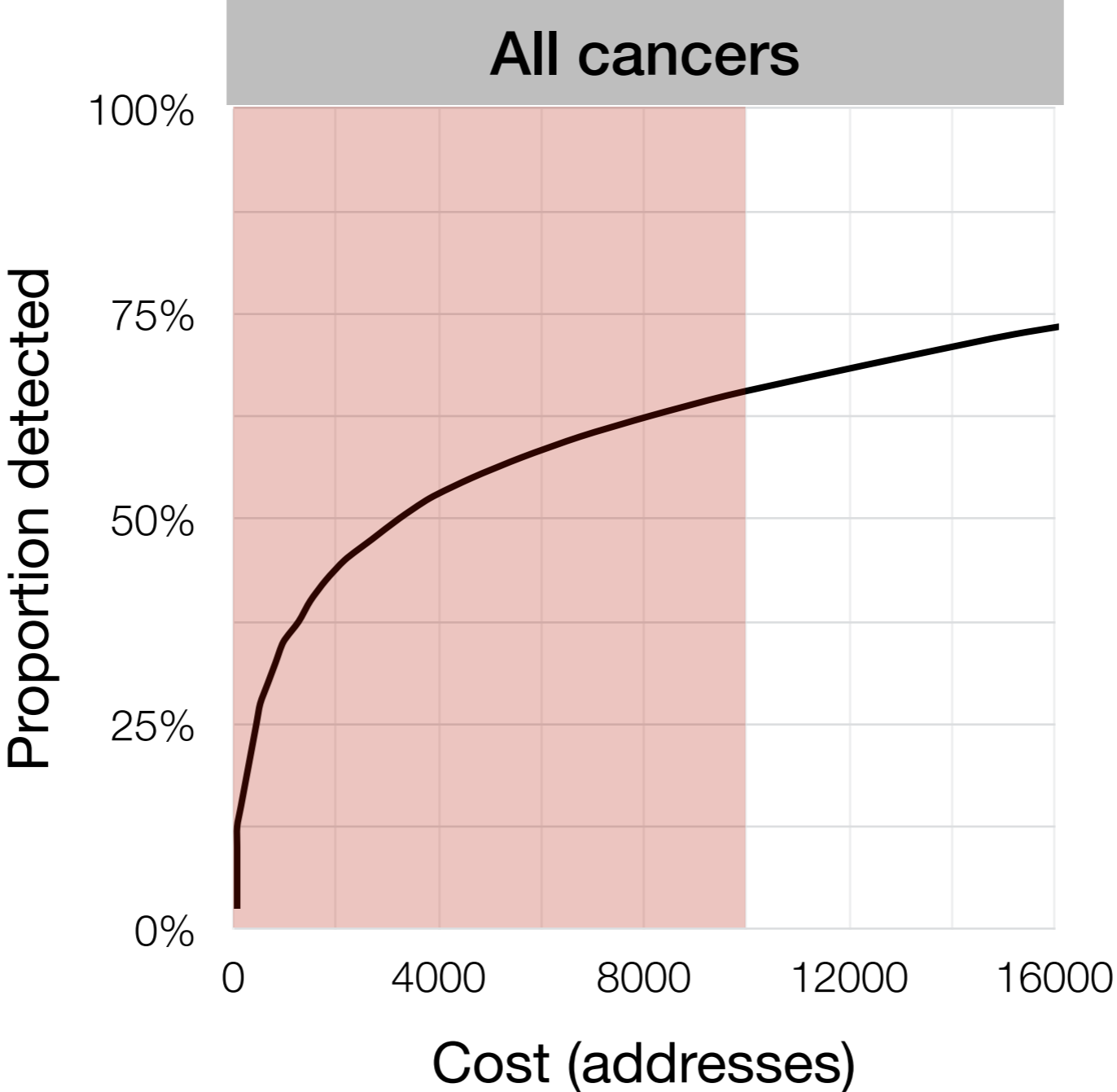
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	1	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
1	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0

Liver
Stomach
Stomach
Pancreas
Pancreas
Lung
Breast
Lung
Lung
Lung
Liver
Breast
Stomach
Pancreas

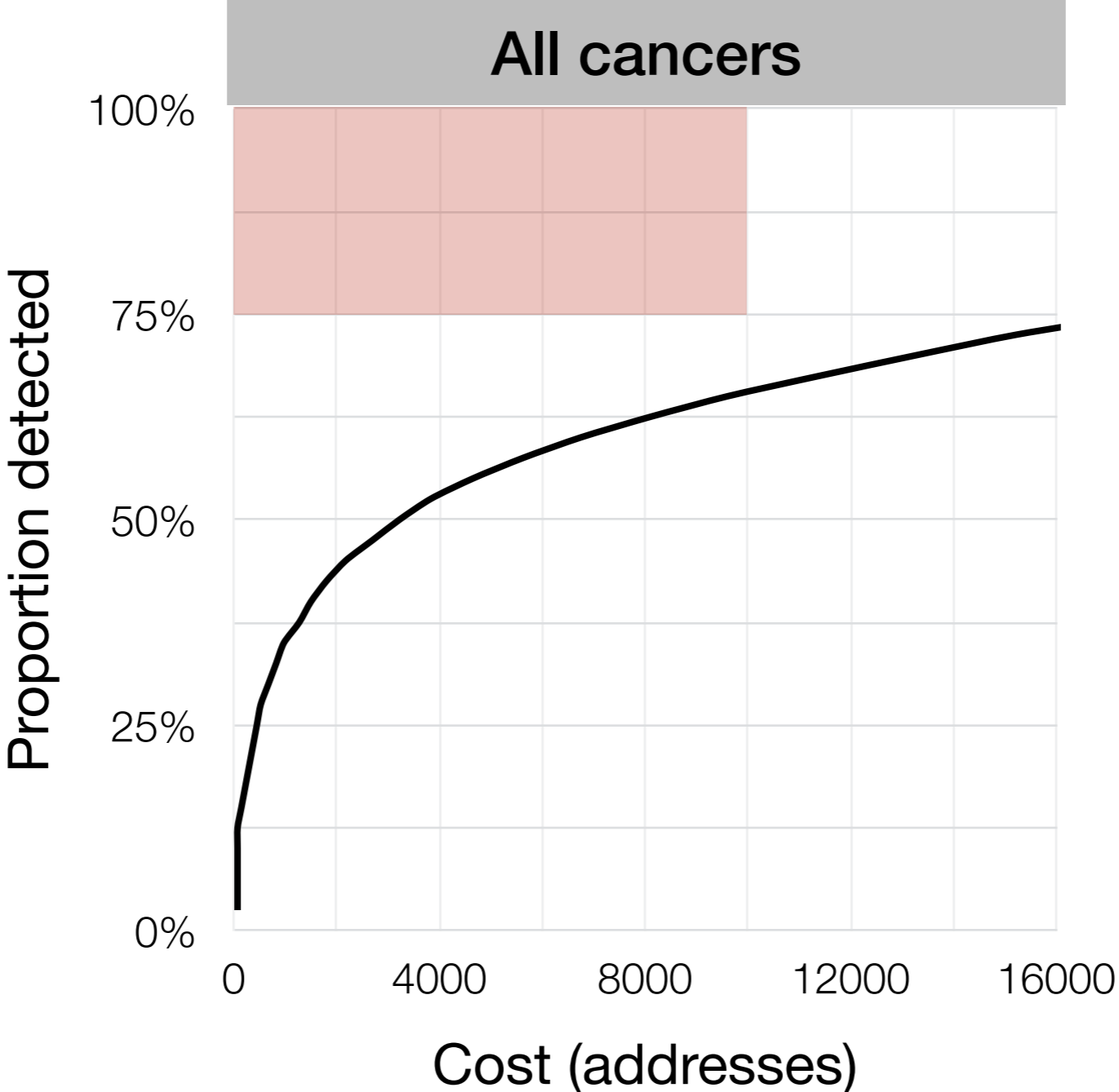
DNA Panels from [Cohen et al.]



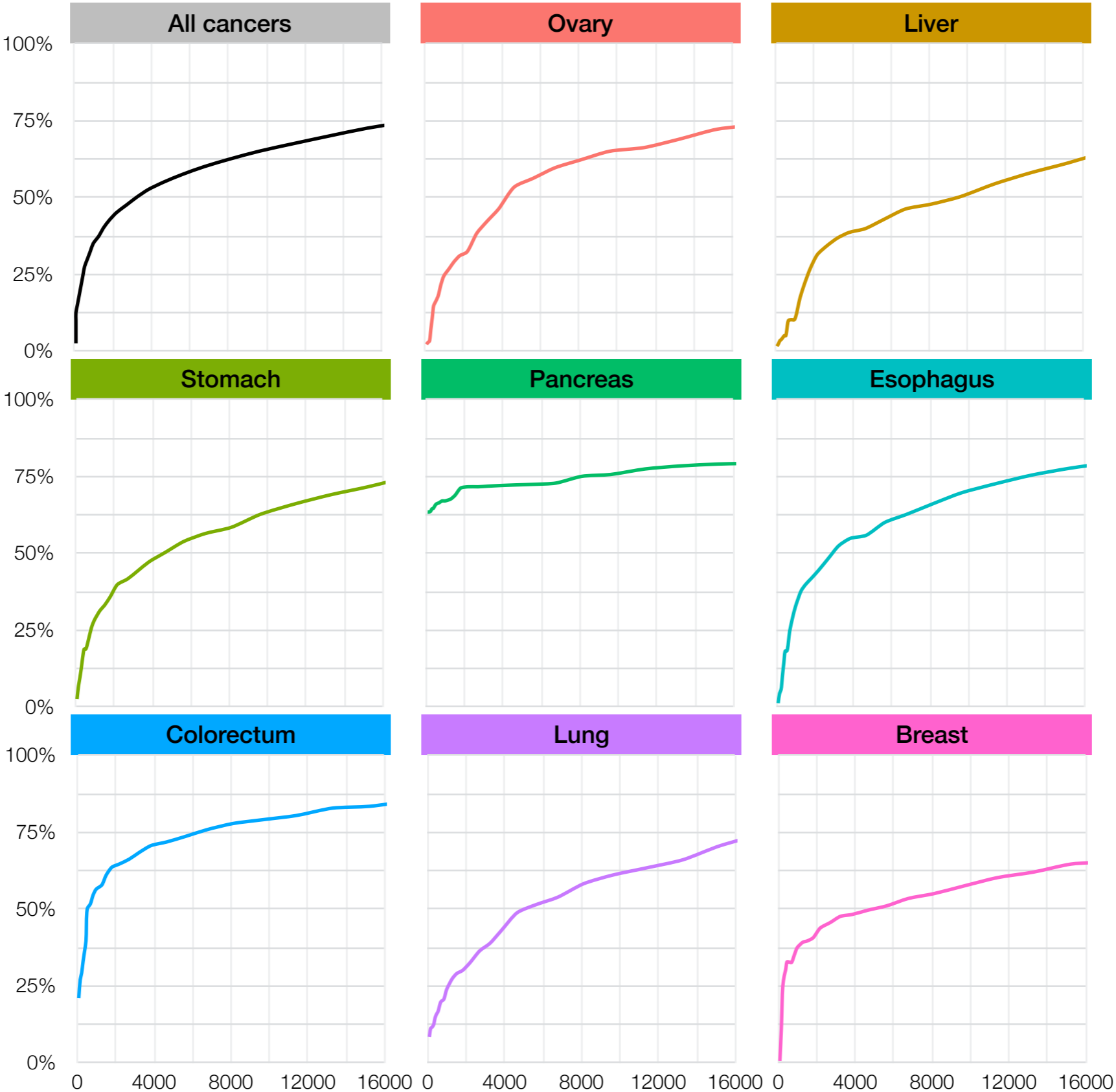
DNA Panels from [Cohen et al.]



DNA Panels from [Cohen et al.]

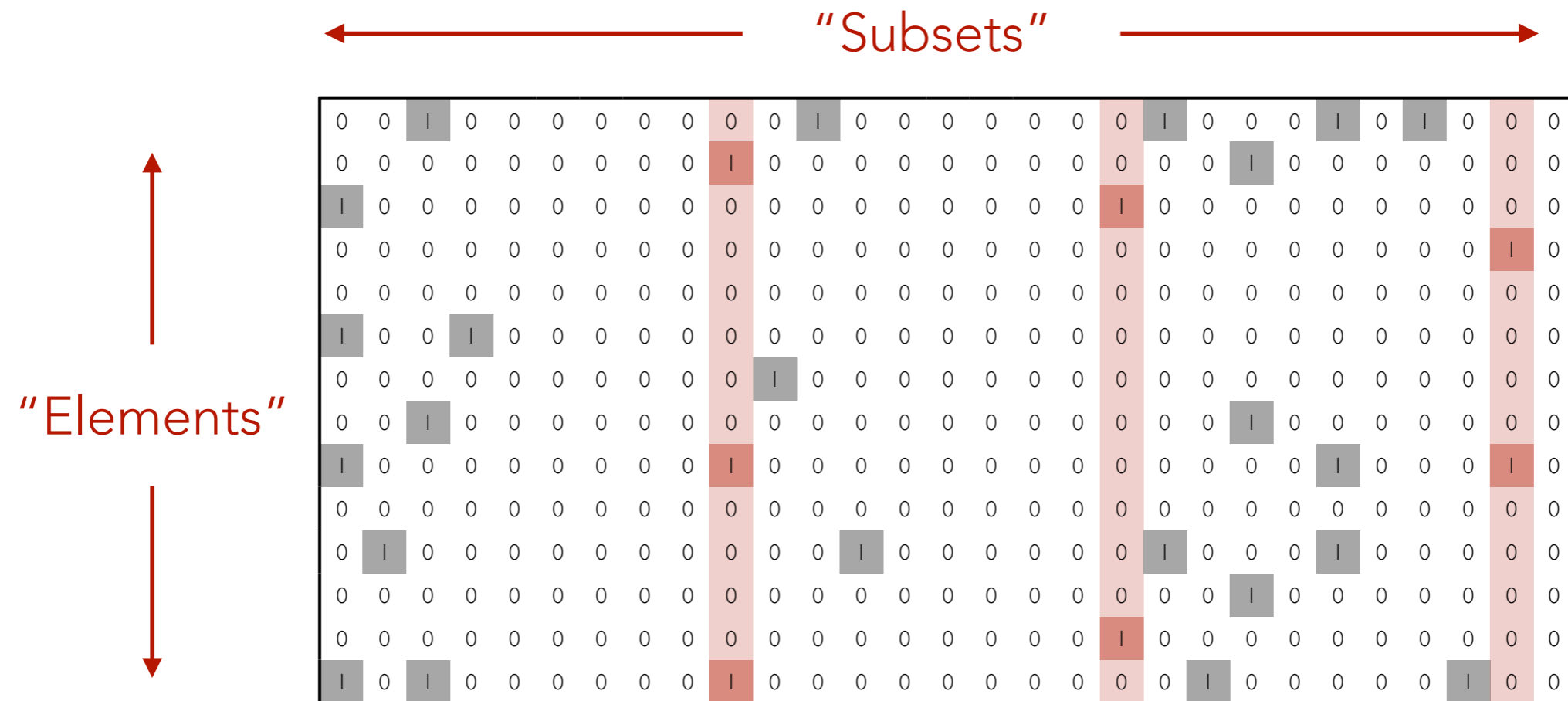


DNA Panels from [Cohen et al.]



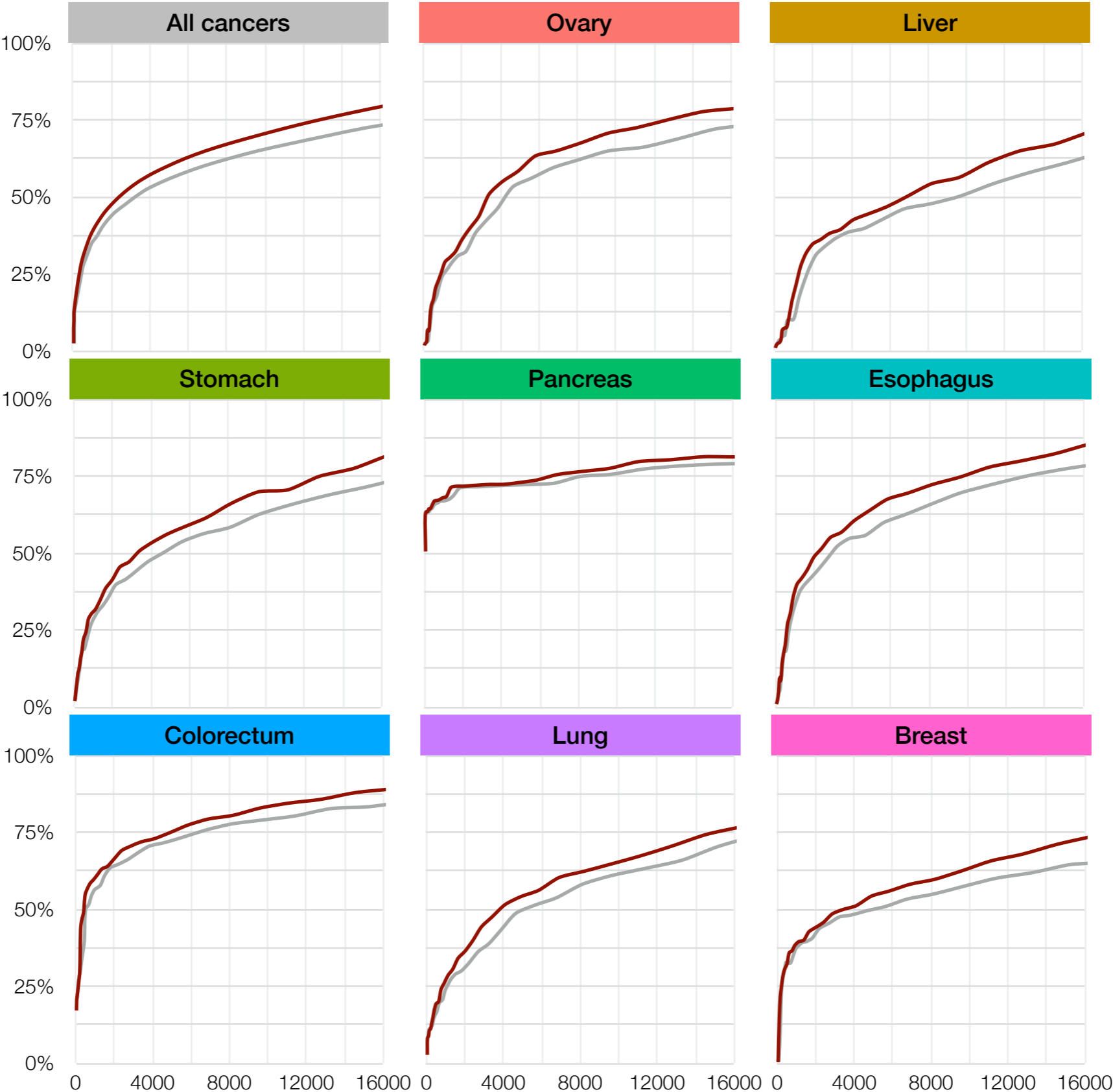
This is an optimization problem!

- *Max Cover*, to be specific

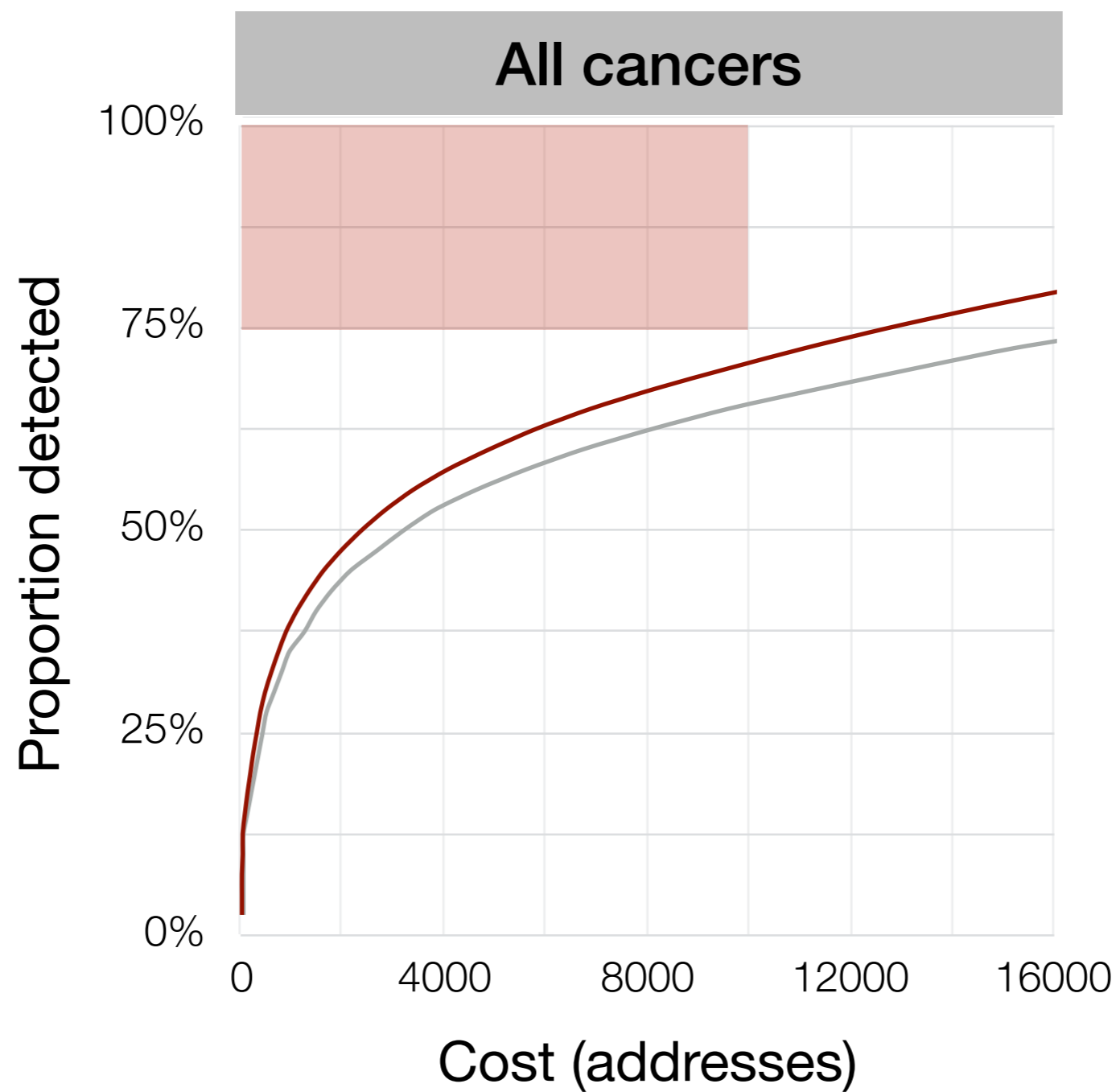


- Huge instance (10^4 elements, 10^9 subsets)
- Still solves in Gurobi

Optimal Panels



Optimal Panels



Adaptive Panels

- **Adaptivity:** perform the test in stages
- Opportunity: the \$100 constraint only needs to hold in **expectation**
- Higher accuracy at same **average** cost

Adaptive Panels

- **Adaptivity:** perform the test in stages
- Opportunity: the \$100 constraint only needs to hold in **expectation**
 - Higher accuracy at same **average** cost
- Challenge: this problem is **hard** (theoretically)

Adaptive Panels

- **Adaptivity:** perform the test in stages
- Opportunity: the \$100 constraint only needs to hold in **expectation**
 - Higher accuracy at same **average** cost
- Challenge: this problem is **hard** (theoretically) and **hard** (practically)

A Simple Model for DNA Mutations

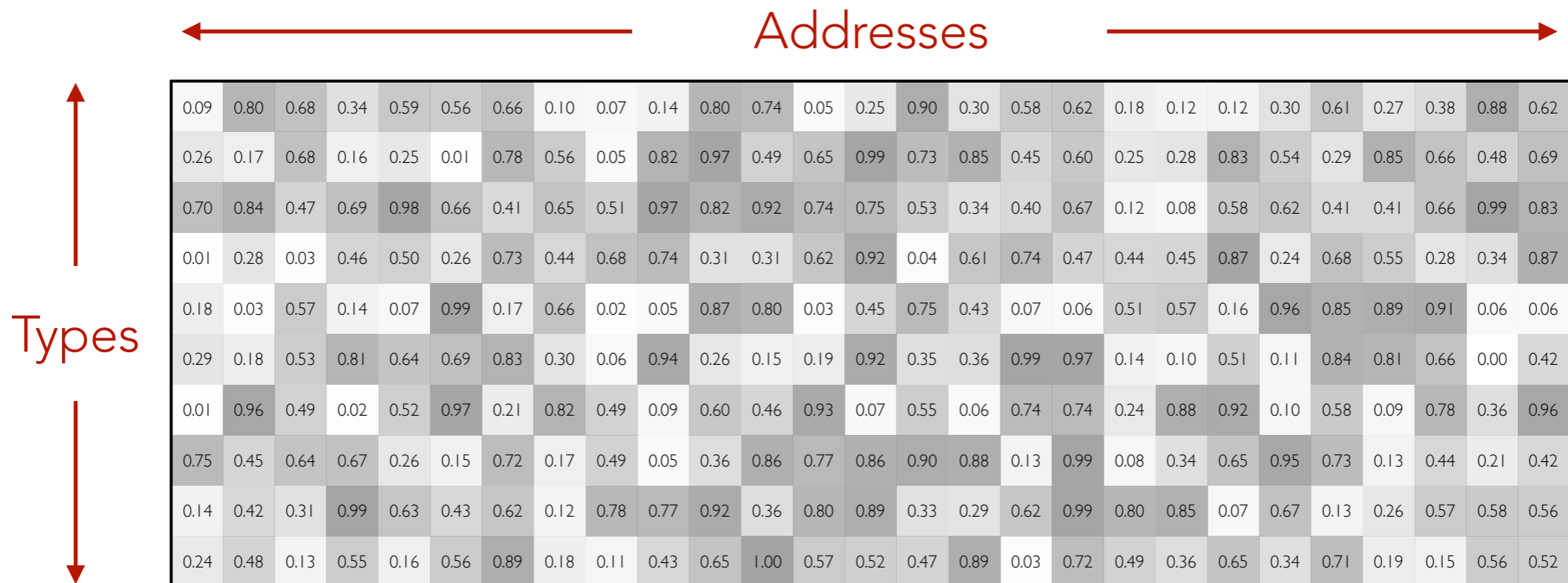
- Cancer *types* $t = 1, \dots, T$

A Simple Model for DNA Mutations

- Cancer **types** $t = 1, \dots, T$
- DNA **addresses** $a = 1, \dots, A$

A Simple Model for DNA Mutations

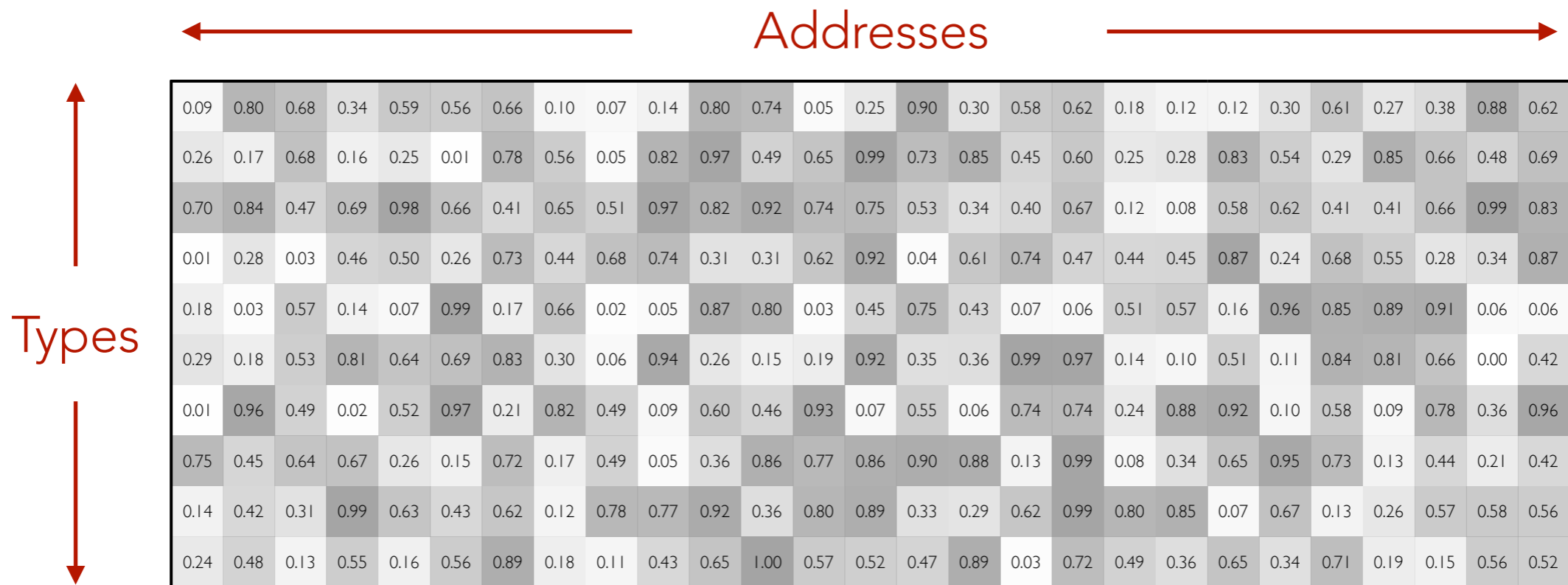
- Cancer **types** $t = 1, \dots, T$
- DNA **addresses** $a = 1, \dots, A$



$$P \in \mathbb{R}^{T \times A}$$

A Simple Model for DNA Mutations

- Cancer **types** $t = 1, \dots, T$
- DNA **addresses** $a = 1, \dots, A$



$$P \in \mathbb{R}^{T \times A}$$

- Sequencing address a on individual of type t :
 - Yields observation $\sim \text{Ber}(P_{ta})$

Problem: Active Sequential Hypothesis Testing

- Unknown cancer type drawn according to (known) **prior**

Problem: Active Sequential Hypothesis Testing

- Unknown cancer type drawn according to (known) **prior**
- **Partial Adaptivity:** select a **sequence** of addresses
 - Same sequence always used
 - Stop anytime, and “guess” the cancer type

Problem: Active Sequential Hypothesis Testing

- Unknown cancer type drawn according to (known) **prior**
- **Partial Adaptivity:** select a **sequence** of addresses
 - Same sequence always used
 - Stop anytime, and “guess” the cancer type
- Constraint: correctly identify type with probability at least $1 - \delta$
- Objective: minimize **expected cost (number of addresses used)**

Theoretical Guarantee

- Recall:
 - A = number of addresses (3 billion)
 - T = number of cancer types (10-100)

Theoretical Guarantee

- Recall:
 - A = number of addresses (3 billion)
 - T = number of cancer types (10-100)
- L = maximum length of test sequence (can think of this as $\sim 10,000$)

Theoretical Guarantee

- Recall:
 - A = number of addresses (3 billion)
 - T = number of cancer types (10-100)
- L = maximum length of test sequence (can think of this as $\sim 10,000$)

	Runtime	Cost Guarantee
Brute Force	A^L	OPT

Theoretical Guarantee

- Recall:
 - A = number of addresses (3 billion)
 - T = number of cancer types (10-100)
- L = maximum length of test sequence (can think of this as $\sim 10,000$)

	Runtime	Cost Guarantee
Brute Force	A^L	OPT
LP Heuristic [Naghshvar,Javidi'13]	$O(AT^2)$	None

Theoretical Guarantee

- Recall:
 - A = number of addresses (3 billion)
 - T = number of cancer types (10-100)
- L = maximum length of test sequence (can think of this as $\sim 10,000$)

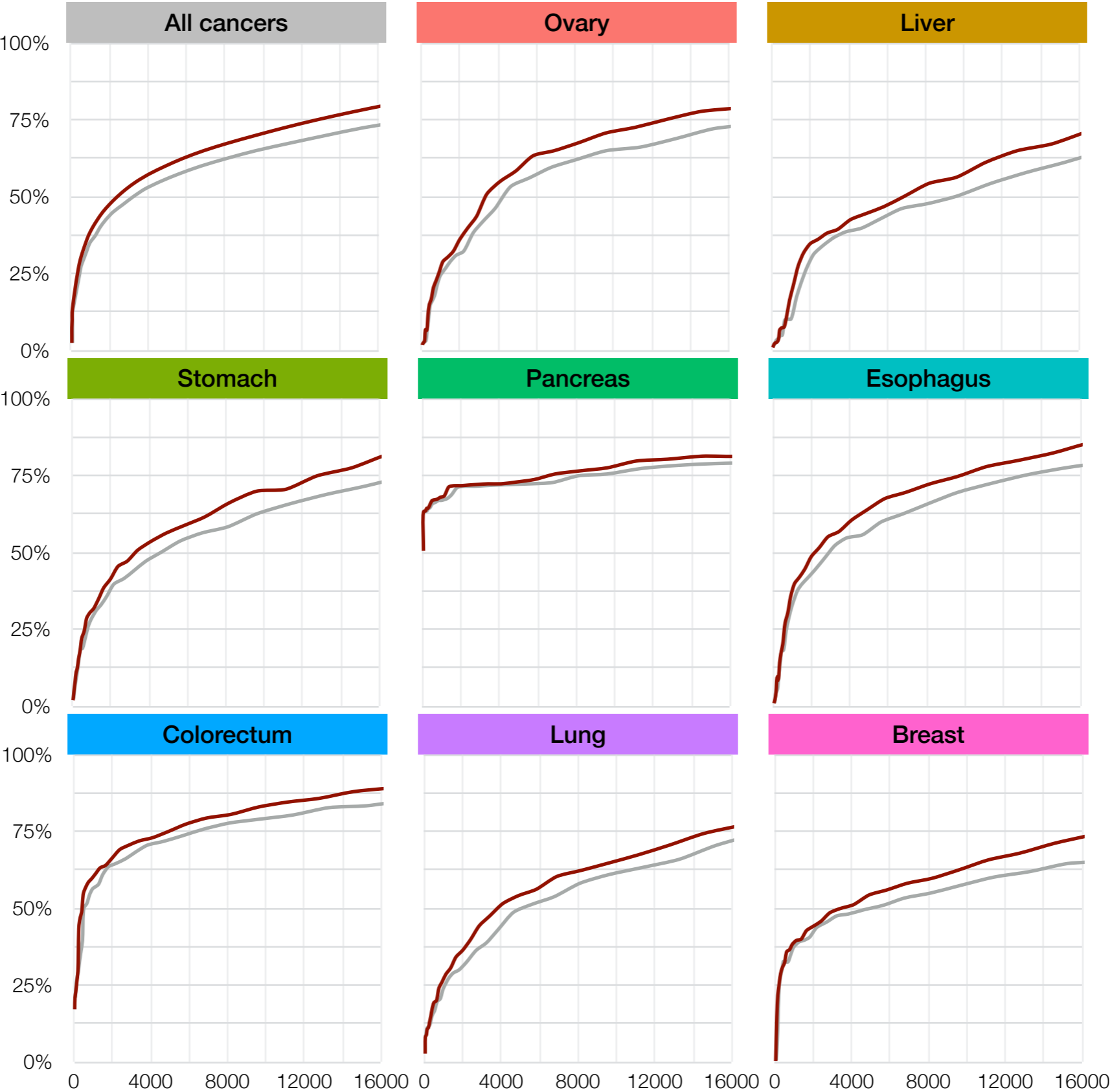
	Runtime	Cost Guarantee
Brute Force	A^L	OPT
LP Heuristic [Naghshvar,Javidi'13]	$O(AT^2)$	None
Our Algorithm	$O(ATL)$	$O(\log T)OPT$

Theoretical Guarantee

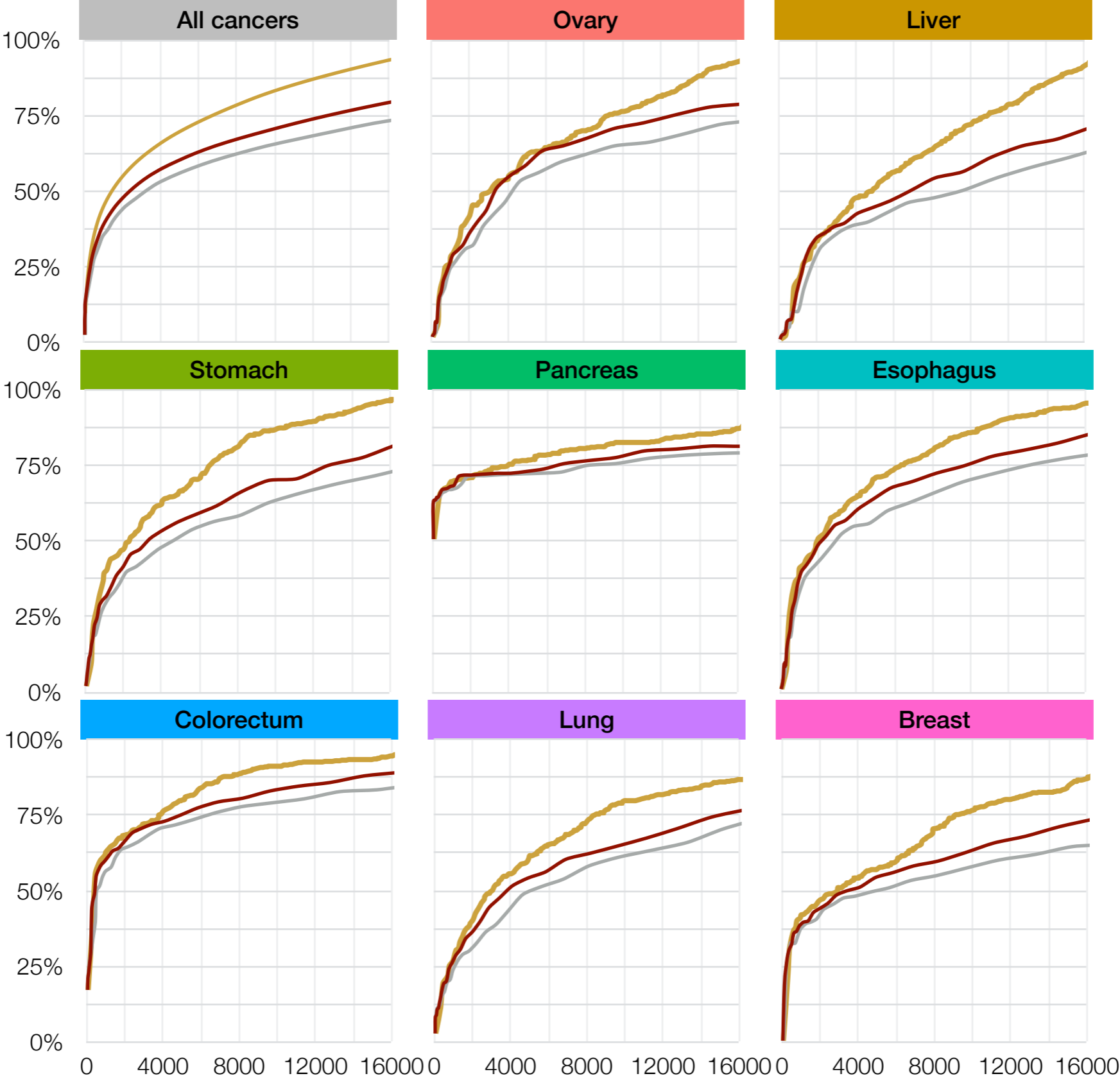
- Recall:
 - A = number of addresses (3 billion)
 - T = number of cancer types (10-100)
- L = maximum length of test sequence (can think of this as $\sim 10,000$)

	Runtime	Cost Guarantee
Brute Force	A^L	OPT
LP Heuristic [Naghshvar,Javidi'13]	$O(AT^2)$	None
Our Algorithm	$O(ATL)$	$O(\log T + \log \log \delta^{-1})OPT$

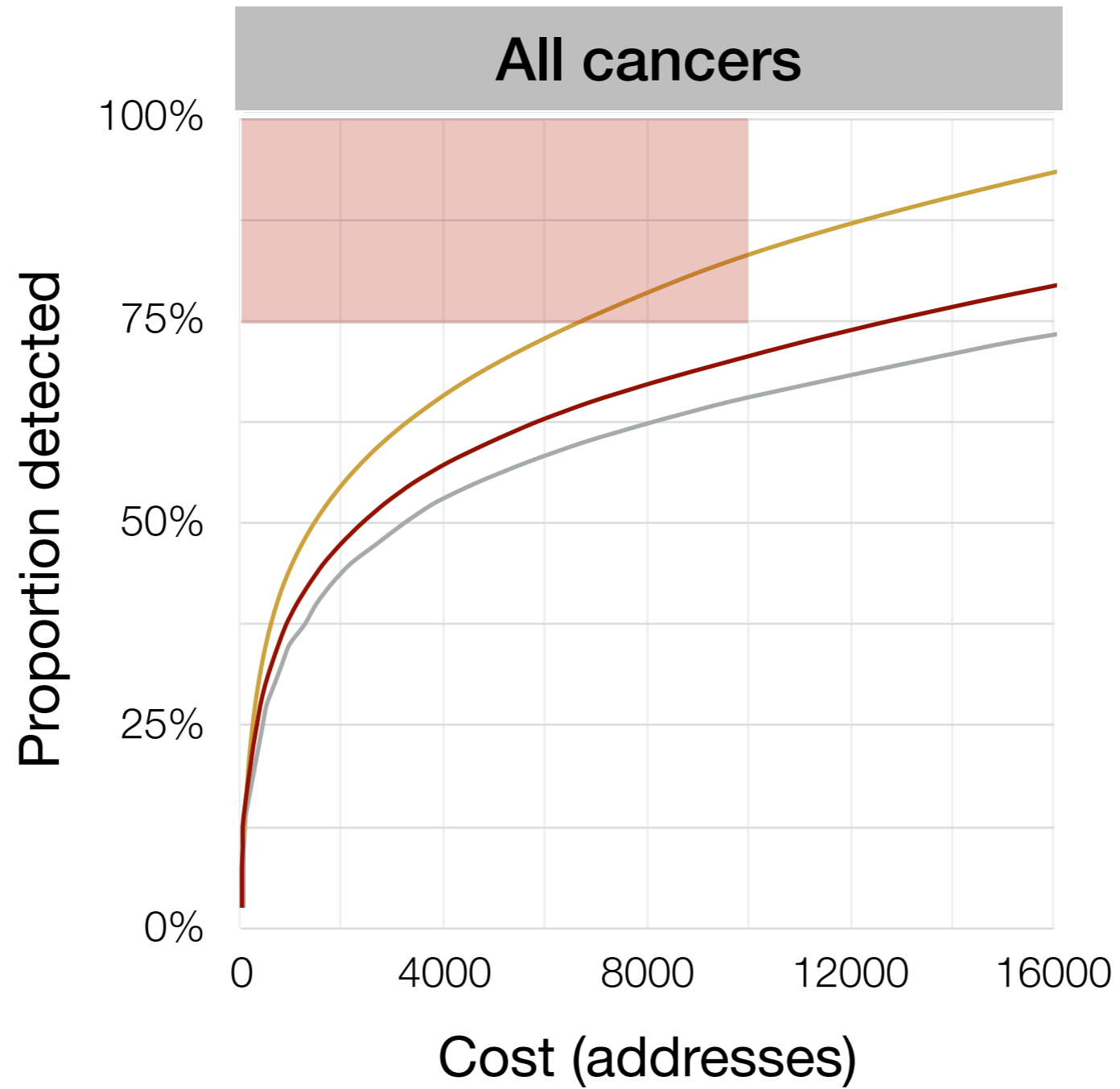
Optimal (Non-adaptive) Panels



(Sub-optimal) Adaptive Panels



(Sub-optimal) Adaptive Panels



Thanks!