Introduction	

Improved FPT algorithm 000000

Induced Haplotyping

Extended Islands of Tractability for Parsimony Haplotyping

Rudolf Fleischer¹, Jiong Guo², Rolf Niedermeier³, Johannes Uhlmann³, Yihui Wang¹, <u>Mathias Weller³</u>, and Xi Wu¹

¹Fudan University Shanghai, ²Universität des Saarlandes, and ³Friedrich-Schiller-Universität Jena

June 22, 2010

 Introduction
 Improved FPT algorithm
 Induced Haplotyping
 Conclusion

 •0000000
 First, Some Biology...
 Conclusion
 Conclusion
 Conclusion



approx. 0.1% of human nucleotide sites differ between individuals

Introduction •0000000 Improved FPT algorithm

Induced Haplotyping

Conclusion 00

First, Some Biology...



approx. 0.1% of human nucleotide sites differ between individuals

Introduction •0000000 Improved FPT algorithm

Induced Haplotyping

Conclusion 00

First, Some Biology...



the sequence of SNPs is called a haplotype

 Introduction
 Improved FPT algorithm
 Induced Haplotyping
 Conclusion

 •ooooooo
 First, Some Biology...
 Conclusion
 Conclusion
 Conclusion



the sequence of SNPs is called a haplotype

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
0●000000	000000	0000	00
First, Some Bi	ology		



humans are diploid \rightarrow 2 chromosome sets haplotype = SNP sequence in one chromosome set genotype = SNP sequence in the combined chromosome sets Introduction 0000000 Improved FPT algorithm

Induced Haplotyping

Conclusion 00

First, Some Biology...



humans are diploid \sim 2 chromosome sets haplotype = SNP sequence in one chromosome set genotype = SNP sequence in the combined chromosome sets

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
0000000	000000	0000	
Motivation			

Haplotype Inference

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
0000000	000000	0000	00
Motivation			

Haplotype Inference

Problems

- impossible to infer haplotypes of just 1 genotype
 → sequence and infer groups/populations
- which explanation should be preferred if there are multiple?
 → parsimony ("Ockham's razor")
- how to perform the actual computation fast?

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
0000000	000000	0000	00
Motivation			

Haplotype Inference

Problems

- impossible to infer haplotypes of just 1 genotype
 → sequence and infer groups/populations
- which explanation should be preferred if there are multiple?
 → parsimony ("Ockham's razor")
- how to perform the actual computation fast?

Preliminary De	finitions		
0000000	000000	0000	00
Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion

Definition (Haplotype, Genotype, Resolve)

```
haplotype h = \text{string over } \{0, 1\}
genotype g = \text{string over } \{0, 1, 2\}
h_1, h_2 resolve g \Leftrightarrow
for all i \in \mathbb{N}, g[i] = h_1[i] = h_2[i] or g[i] = 2 and h_1[i] \neq h_2[i]
multiset H \rightsquigarrow \text{res}(H)
multiset H resolves multiset G \Leftrightarrow G \subseteq \text{res}(H)
```

Example											
haplotype1:	0	0	1	0	1	1	1	0	0	0	1
haplotype2:	0	0	1	1	0	1	0	0	1	1	1
genotype:	0	0	1	2	2	1	2	0	2	2	1

Preliminary De	finitions		
0000000	000000	0000	00
Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion

Definition (Haplotype, Genotype, Resolve)

```
haplotype h = \text{string over } \{0, 1\}
genotype g = \text{string over } \{0, 1, 2\}
h_1, h_2 resolve g \Leftrightarrow
for all i \in \mathbb{N}, g[i] = h_1[i] = h_2[i] or g[i] = 2 and h_1[i] \neq h_2[i]
multiset H \rightsquigarrow \text{res}(H)
multiset H resolves multiset G \Leftrightarrow G \subseteq \text{res}(H)
```

Example											
haplotype1:	0	0	1	0	1	1	1	0	0	0	1
haplotype2:	0	0	1	1	0	1	0	0	1	1	1
genotype:	0	0	1	2	2	1	2	0	2	2	1

Preliminary [Definitions		
0000000	000000	0000	00
Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion

Definition (Haplotype, Genotype, Resolve)

```
haplotype h = \text{string over } \{0, 1\}
genotype g = \text{string over } \{0, 1, 2\}
h_1, h_2 resolve g \Leftrightarrow
for all i \in \mathbb{N}, g[i] = h_1[i] = h_2[i] or g[i] = 2 and h_1[i] \neq h_2[i]
multiset H \rightsquigarrow \text{res}(H)
multiset H resolves multiset G \Leftrightarrow G \subseteq \text{res}(H)
```

Example											
haplotype1:	0	0	1	1	1	1	1	0	0	0	1
haplotype2:	0	0	1	0	0	1	0	0	1	1	1
genotype:	0	0	1	2	2	1	2	0	2	2	1

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	000000	0000	
Preliminary De	finitions		

Definition (Haplotype Graph)

Given H and $G \subseteq res(H)$, define the haplotype graph of H and G:

- |H| vertices (labeled by H)
- |G| edges (labeled by G)

• haplotypes of the endpoints of an edge resolve its genotype



Introduction	
000000000	

Improved FPT algorithm

Induced Haplotyping

Conclusion

Haplotype Inference by Parsimony

Definition (Haplotype Inference by Parsimony)

Input: multiset *G* of length-*m* genotypes, integer $k \ge 0$ **Question**: \exists multiset *H* of *k* haplotypes that resolves *G* ?

Example haplotype graph haplotypes genotypes

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000			
Previous Work			

Complexity

- NP-hard [Halldórsson et al., DMTCS '03]
- APX-hard [Lancia et al., INFORMS '04]
- many special cases in P (e.g. [Lancia & Rizzi, ORL '06])

Algorithms

- Factor-2^{d-1}-Approximation [Lancia & Rizzi, ORL '06]
- $O(2^{|G| \cdot d})$ Branch&Bound [Wang & Xu, Bioinformatics '03]
- $O(k^{2k^2} \cdot m)$ Algorithm [Sharan et al., TCBB '06]

k := #haplotypes, m := genotype length, $d := \max \# 2$'s in a genotype

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
○○○○○○○	000000	0000	00
Our Contribut	ion		

- $k^{4k} \cdot \text{poly}(|G|, m)$ time algorithm
- polynomial-time solvable special case: "Induced Haplotype Inference"

00000000		0000	00		
Inference Graph					

Definition (Inference Graph)

inference graph Γ of G = an order-k graph with edges **consistently** labeled by the genotypes in G



Introduction 00000000 Improved FPT algorithm 00000

Induced Haplotyping

Conclusion 00

If Only We Had an Inference Graph...

Observation (non-bipartite components of Γ)

 $C = cycle \text{ in } \Gamma$ \rightsquigarrow for all i, $|\{g \in C \mid g[i] = 2\}|$ is even non-bipartite components of $\Gamma \rightsquigarrow O(|\Gamma| \cdot m)$ time



Introduction 00000000 Improved FPT algorithm

Induced Haplotyping

Conclusion 00

If Only We Had an Inference Graph...

Observation (bipartite components of Γ) all g[i] = 2 for some $i \Rightarrow$ choose arbitrarily bipartite components of $\Gamma \rightarrow O(|\Gamma| \cdot m)$ time



Introduction Improved FPT algorithm Induced Haplotyping Conclusion oc

If Only We Had an Inference Graph...

Observation

(G, k) yes-instance with solution H $\Rightarrow \exists \Gamma$ extendable $(O(|\Gamma| \cdot m) \text{ time})$ to a haplotype graph of H and G Introduction Improved FPT algorithm Induced Haplotyping conclusion oc

If Only We Had an Inference Graph...

Observation

(G, k) yes-instance with solution H

 $\Rightarrow \exists \ \Gamma \ extendable \ (O(|\Gamma| \cdot m) \ time)$ to a haplotype graph of H and G

algorithmic idea: guess Г

 Introduction
 Improved FPT algorithm
 Induced Haplotyping
 Conclusion

 0000000
 0000000
 0000
 0000
 0000

If Only We Had an Inference Graph...

Observation

(G, k) yes-instance with solution H

 $\Rightarrow \exists \ \Gamma \ extendable \ (O(|\Gamma| \cdot m) \ time)$ to a haplotype graph of H and G

algorithmic idea: guess Γ **better idea:** guess a "spanning" subgraph of Γ

Algorithm	Outline		
Introduction 00000000	Improved FPT algorithm ○○○○●○	Induced Haplotyping	Conclusion

Algorithm

- **Ο** guess "spanning" subgraph of Γ
- **2** infer the haplotype multiset $H \rightsquigarrow O(k \cdot m)$ time
- So check whether H resolves $G \rightsquigarrow O(k^2 \cdot m)$ time

Introduction	

Improved FPT algorithm

Induced Haplotyping

Conclusion 00

Algorithm Outline

Algorithm

guess "spanning" subgraph of F

- guess a size-k genotype subset of $G \rightsquigarrow O(k^{2k})$ possibilities
- **②** for these genotypes, guess 2 (of k) vertices $\rightsquigarrow O(k^{2k})$ possibilities
- **2** infer the haplotype multiset $H \rightsquigarrow O(k \cdot m)$ time
- So check whether H resolves $G \rightsquigarrow O(k^2 \cdot m)$ time

Introduction	

Improved FPT algorithm

Induced Haplotyping

Conclusion 00

Algorithm Outline

Algorithm

guess "spanning" subgraph of F

- guess a size-k genotype subset of $G \rightsquigarrow O(k^{2k})$ possibilities
- **②** for these genotypes, guess 2 (of k) vertices $\rightsquigarrow O(k^{2k})$ possibilities
- ② infer the haplotype multiset $H \rightsquigarrow O(k \cdot m)$ time
- So check whether H resolves $G \rightsquigarrow O(k^2 \cdot m)$ time

Theorem

Haplotype Inference by Parsimony can be solved in $O(k^{4k+2} \cdot m)$ time.

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	○○○○○●		00
Example			



Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	○○○○○●	0000	00
Example			



Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	○○○○○●	0000	00
Example			



Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000		0000	00
Example			



Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000		0000	00
Example			



Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	○○○○○●		00
Example			



00000000	000000	•000	00
Induced Haplo	type Inference		

recall: *H* resolves $G \Leftrightarrow G \subseteq \text{res}(H)$ what if G = res(H)? \rightsquigarrow haplotype graph is a clique

Definition

 $G = \operatorname{res}(H) \rightsquigarrow H \text{ induces } G$

0000000		0000	00
Induced Ha	aplotype Inference		

recall: *H* resolves $G \Leftrightarrow G \subseteq \text{res}(H)$ what if G = res(H)? \rightsquigarrow haplotype graph is a clique

Definition

 $G = \operatorname{res}(H) \rightsquigarrow H \text{ induces } G$

Definition (Induced Haplotype Inference by Parsimony)

Input: multiset *G* of length-*m* genotypes **Question**: \exists multiset *H* of haplotypes that induces *G* ?

Observatio	ns & Algorithm		
00000000			00
Introduction	Improved EPT algorithm	Induced Hanletyning	Conclusion

G can be partitioned "nicely" into G_0 , G_1 , and G_2 .

Obconvoti	one le Algorithm		00
00000000	000000	000	00

G can be partitioned "nicely" into G_0 , G_1 , and G_2 .

Observation

$${\it G}_2 \neq \emptyset$$
 but ${\it G}_0 = \emptyset$ or ${\it G}_1 = \emptyset ~~ \rightsquigarrow ~ poly$

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
0000000	000000	○●○○	
Observations &	& Algorithm		

G can be partitioned "nicely" into G_0 , G_1 , and G_2 .

Observation

$$G_2 \neq \emptyset \text{ but } G_0 = \emptyset \text{ or } G_1 = \emptyset \ \rightsquigarrow \text{ poly}$$

 $|G_0| = |G_1| = 1 \quad \rightsquigarrow \text{ poly (although we may get 2 solutions)}$

Strategy: Divide & Conquer

divide-step base cases

merge-step

00000000		0000	00
Observations	& Algorithm		

G can be partitioned "nicely" into G_0 , G_1 , and G_2 .

Observation

$$G_2 \neq \emptyset \text{ but } G_0 = \emptyset \text{ or } G_1 = \emptyset \ \rightsquigarrow \text{ poly}$$

Strategy: Divide & Conquer	
divide-step → OK	
base cases	
merge-step	

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
0000000	000000	○●○○	
Observations &	& Algorithm		

G can be partitioned "nicely" into G_0 , G_1 , and G_2 .

Observation

$${\it G}_2 \neq \emptyset \text{ but } {\it G}_0 = \emptyset \text{ or } {\it G}_1 = \emptyset \ \rightsquigarrow \text{ poly}$$

Strategy: Divide & Conquer	
divide-step \rightsquigarrow OK	
base cases $\rightsquigarrow OK$	
merge-step	

00000000		0000	00
Observations	& Algorithm		

G can be partitioned "nicely" into G_0 , G_1 , and G_2 .

Observation

$${\it G}_2 \neq \emptyset \text{ but } {\it G}_0 = \emptyset \text{ or } {\it G}_1 = \emptyset \ \rightsquigarrow \text{ poly}$$

Strategy: Div	ide & Conquer	
divide-step	$\sim OK$	
base cases	$\sim OK$	
merge-step	\rightsquigarrow problem!	

Observation	s & Algorithm		
0000000	000000	0000	00
Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion

G can be partitioned "nicely" into G_0 , G_1 , and G_2 .

Observation

$${\it G}_2 \neq \emptyset \text{ but } {\it G}_0 = \emptyset \text{ or } {\it G}_1 = \emptyset \ \rightsquigarrow \text{ poly}$$

 $|G_0| = |G_1| = 1 \quad \rightsquigarrow \text{ poly (although we may get 2 solutions)}$

Strategy: Div	ide & Conquer	
divide-step	~→ OK	
base cases $\rightarrow OK$		
merge-step	\rightarrow problem!	

need to find a way to compute H_1 for given H_0 and G_2

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	000000	○○●○	00
Extending a S	ubclique Solution		

Let H_0 induce G_0 and let g be a genotype in G_2 with the smallest number of 2's. \rightarrow All $h \in H_0$ that are consistent with g are equal.

Introduction	

Improved FPT algorithm

Induced Haplotyping

Conclusion 00

Extending a Subclique Solution

Observation

Let H_0 induce G_0 and let g be a genotype in G_2 with the smallest number of 2's. \rightarrow All $h \in H_0$ that are consistent with g are equal.



Introduction	

Improved FPT algorithm

Induced Haplotyping

Conclusion 00

Extending a Subclique Solution

Observation

Let H_0 induce G_0 and let g be a genotype in G_2 with the smallest number of 2's. \rightarrow All $h \in H_0$ that are consistent with g are equal.



 Introduction
 Improved FPT algorithm
 Induced Haplotyping
 Conclusion

 Concluding The Induced Problem

Divide & Conquer like	algorithm
divide steps	$O(G \cdot k)$
base solutions	$O(G \cdot m)$
extends (merges)	$O(G \cdot k \cdot m)$
all in all	$O(G \cdot k \cdot m)$

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	000000	0000	●0
Conclusion			

what we saw...

- Induced Haplotyping in $O(|G| \cdot k \cdot m)$ time
- $2^{O(k^2 \log k)}$ -time algorithm improved to $2^{O(k \log k)}$

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	000000	0000	●O
Conclusion			

what we saw...

- Induced Haplotyping in $O(|G| \cdot k \cdot m)$ time
- $2^{O(k^2 \log k)}$ -time algorithm improved to $2^{O(k \log k)}$

also in the paper

- results also hold for sets instead of multisets
- results basically also hold for constrained variant
- data reduction $(O(2^k \cdot k^2)$ -bit kernel)

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
00000000	000000	0000	●O
Conclusion			

what we saw...

- Induced Haplotyping in $O(|G| \cdot k \cdot m)$ time
- $2^{O(k^2 \log k)}$ -time algorithm improved to $2^{O(k \log k)}$

also in the paper

- results also hold for sets instead of multisets
- results basically also hold for constrained variant
- data reduction $(O(2^k \cdot k^2)$ -bit kernel)

future work

- find polynomial kernel (or prove nonexistence)
- distance from triviality measures
- find 2^{O(k)} time algorithm

Introduction	Improved FPT algorithm	Induced Haplotyping	Conclusion
0000000	000000	0000	O•

Thank you