

File Synchronization from Insertions and Deletions

Nicolas Bitouzé

<http://bitouze.bol.ucla.edu>

LORIS Lab, UCLA

Joint work with:

Fred Sala (UCLA),

Lara Dolecek (UCLA),

S. M. Sadegh Tabatabaei Yazdi (Qualcomm)

CoDESS Kickoff Meeting

Sept. 19th, 2013

Problem Setting

File X: h d c e a t g b k u j r v c x f q...

File Y: h d e a k t g v b j r v c r f s q...

Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

Problem Setting

File X: h d c e a t g b k u j r v c x f q...

File Y: h d e a k t g v b j r v c r f s q...

Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

Problem Setting

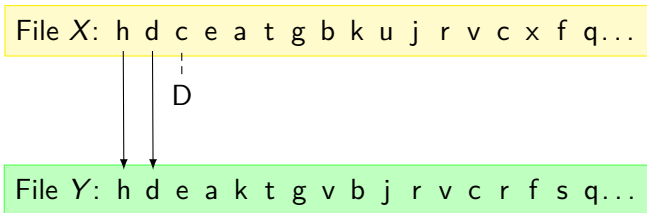
File X: h d c e a t g b k u j r v c x f q...

File Y: h d e a k t g v b j r v c r f s q...

Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

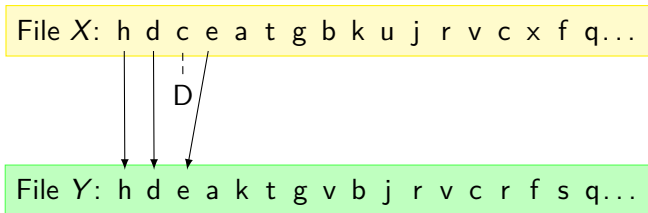
Problem Setting



Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

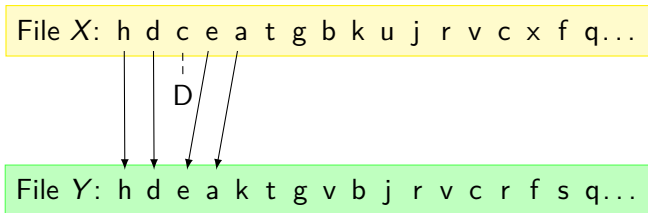
Problem Setting



Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

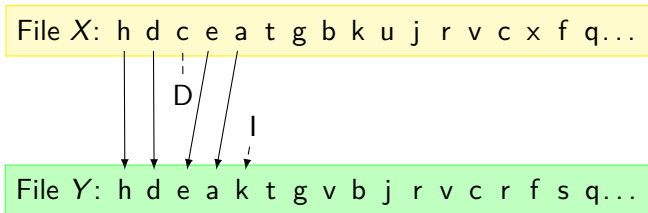
Problem Setting



Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

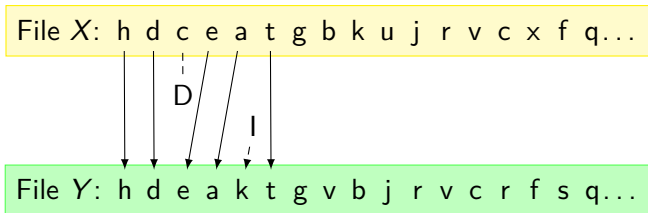
Problem Setting



Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

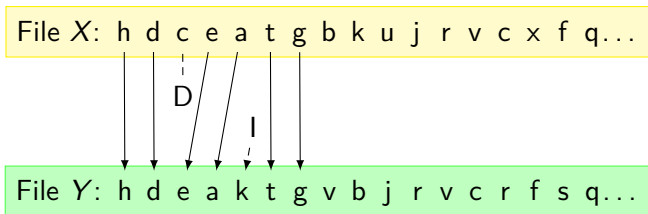
Problem Setting



Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

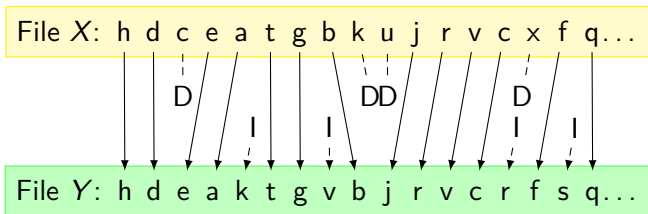
Problem Setting



Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

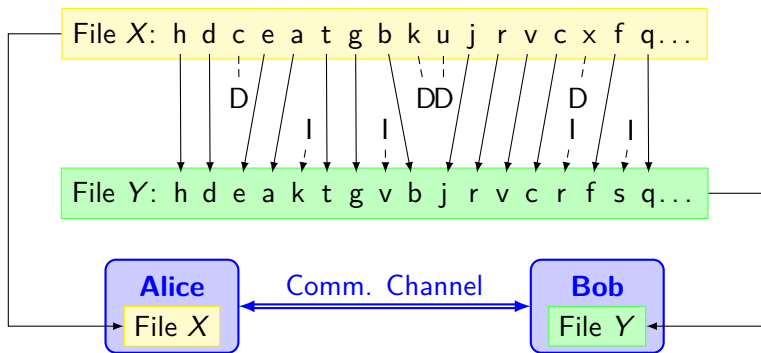
Problem Setting



Small **rate** of edits:

Each symbol is edited with probability $\beta \ll 1$.

Problem Setting




Goal: Interactive Communication Scheme

Allow Bob to recover X from Y :

- with low probability of error,
- with low communication cost.

Prior Related Work

- Scheme that corrects **a single edit** (non-binary):
 -  [G. M. Tenengolts](#), “Nonbinary codes, correcting single deletion or insertion”, 1984.

Prior Related Work

- Scheme that corrects a **single edit** (non-binary):
 - 📄 G. M. Tenengolts, “Nonbinary codes, correcting single deletion or insertion”, 1984.
- Scheme that corrects a **fixed number of edits** (binary):
 - 📄 R. Venkataramanan, H. Zhang, and K. Ramchandran, “Interactive low-complexity codes for synchronization from deletions and insertions”, 2012.

Prior Related Work

- Scheme that corrects a **single edit** (non-binary):
 - 📄 G. M. Tenengolts, “Nonbinary codes, correcting single deletion or insertion”, 1984.
- Scheme that corrects a **fixed number of edits** (binary):
 - 📄 R. Venkataramanan, H. Zhang, and K. Ramchandran, “Interactive low-complexity codes for synchronization from deletions and insertions”, 2012.
- Theoretical bound for the **fixed rate of edits** case:
 - 📄 N. Ma, K. Ramchandaran and D. Tse, “Efficient file synchronization: a distributed source coding approach”, 2011.

Prior Related Work

- Scheme that corrects a **single edit** (non-binary):
 - 📄 G. M. Tenengolts, “Nonbinary codes, correcting single deletion or insertion”, 1984.
- Scheme that corrects a **fixed number of edits** (binary):
 - 📄 R. Venkataramanan, H. Zhang, and K. Ramchandran, “Interactive low-complexity codes for synchronization from deletions and insertions”, 2012.
- Theoretical bound for the **fixed rate of edits** case:
 - 📄 N. Ma, K. Ramchandaran and D. Tse, “Efficient file synchronization: a distributed source coding approach”, 2011.
- Our work: scheme for **fixed rate of edits**:
 - 📄 N. Bitouzé and L. Dolecek, “Synchronization from insertions and deletions under a non-binary, non-uniform source”, 2013.
 - 📄 S. M. S. Tabatabaei and L. Dolecek, “A deterministic, polynomial-time protocol for synchronizing from deletions”, 2013.

Overview of our Synchronization Scheme

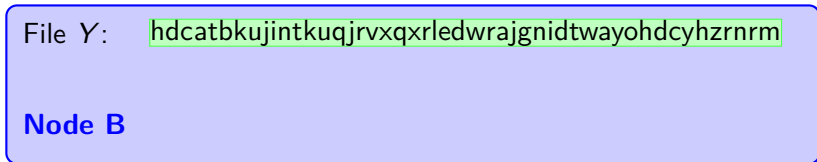
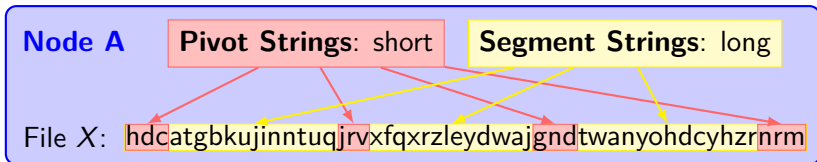
Node A

File X: `hdcatgbkujinntuqjrvxfqxrzleydwajgndtwanyohdcyhzrnm`

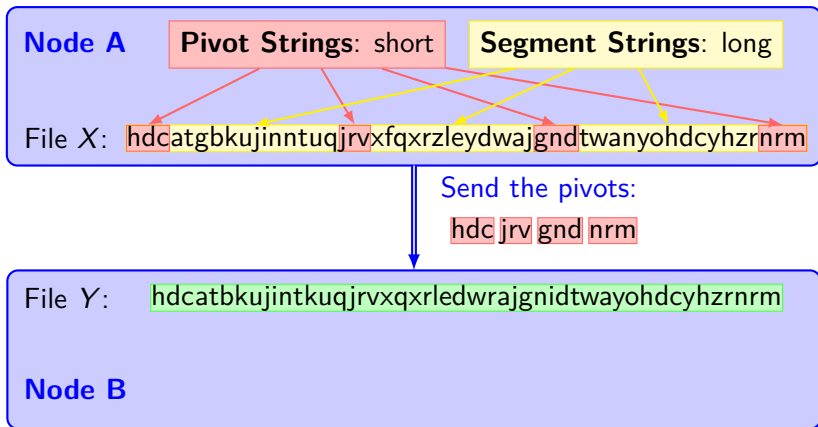
File Y: `hdcatbkujintkuqjrvxqxrledwrajgnidtwayohdcyhzrnm`

Node B

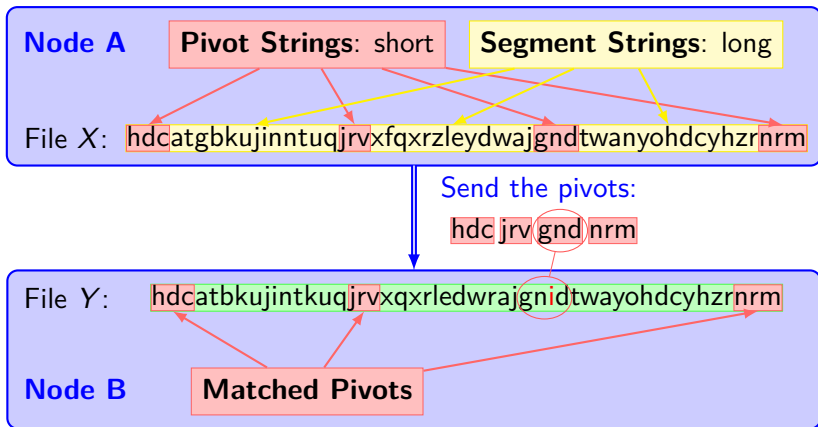
Overview of our Synchronization Scheme



Overview of our Synchronization Scheme

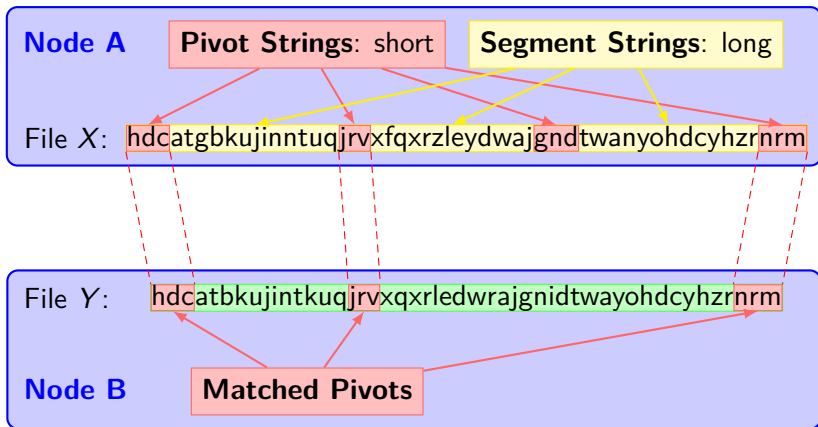


Overview of our Synchronization Scheme



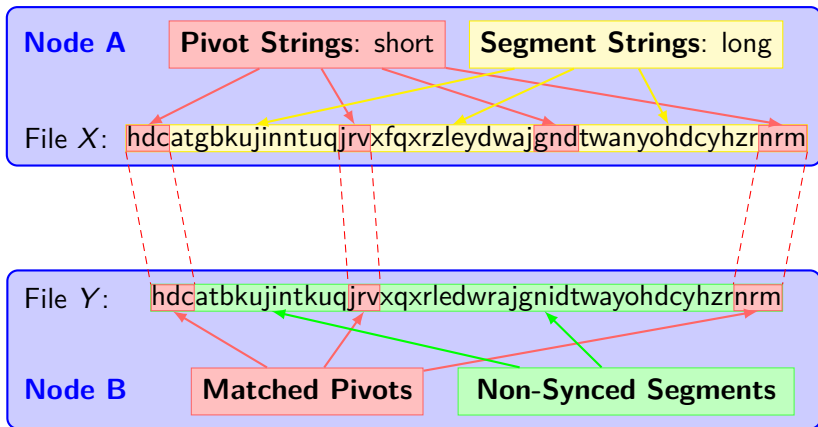
- 1 **Matching Module:** Matches the pivot strings.

Overview of our Synchronization Scheme



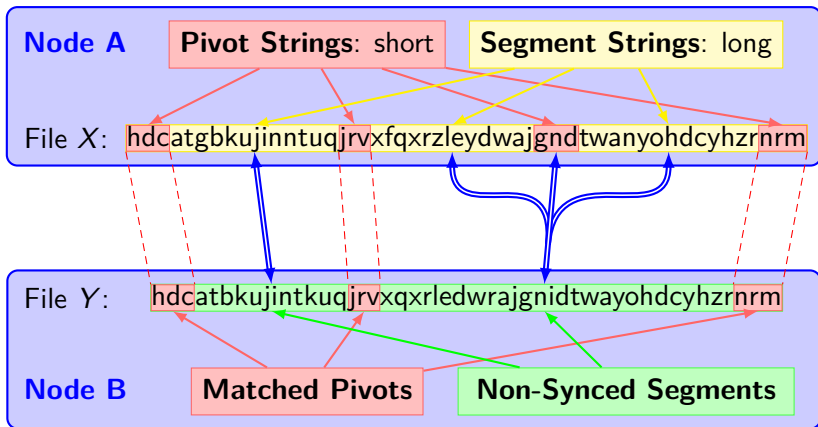
- 1 **Matching Module:** Matches the pivot strings.

Overview of our Synchronization Scheme



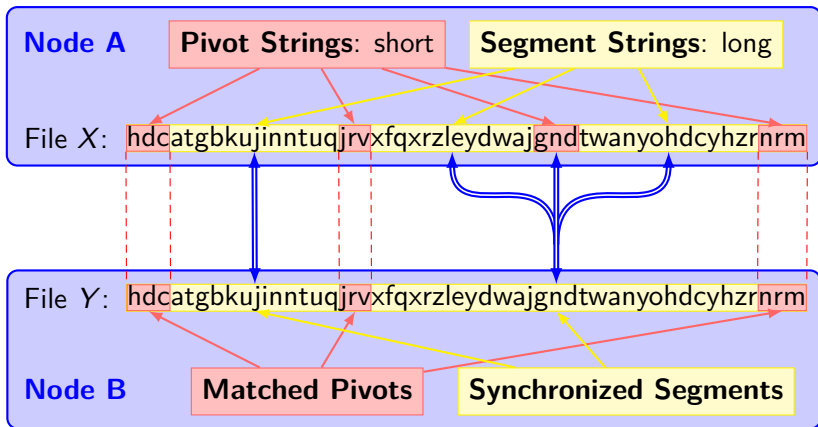
- 1 **Matching Module:** Matches the pivot strings.

Overview of our Synchronization Scheme



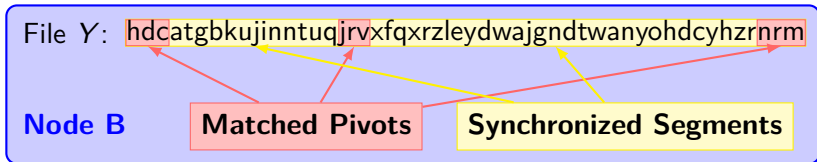
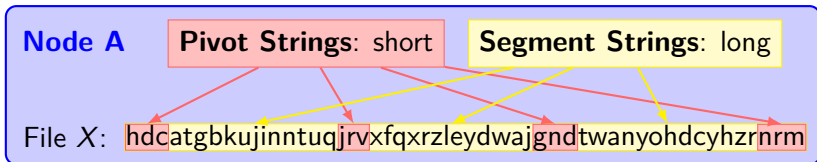
- 1 **Matching Module:** Matches the pivot strings.
- 2 **Edit Recovery Module:** Synchronizes the segment strings.

Overview of our Synchronization Scheme



- 1 **Matching Module:** Matches the pivot strings.
- 2 **Edit Recovery Module:** Synchronizes the segment strings.

Overview of our Synchronization Scheme



- 1 **Matching Module:** Matches the pivot strings.
- 2 **Edit Recovery Module:** Synchronizes the segment strings.
- 3 **Channel Coding Module:** Recovers from residual errors if any.

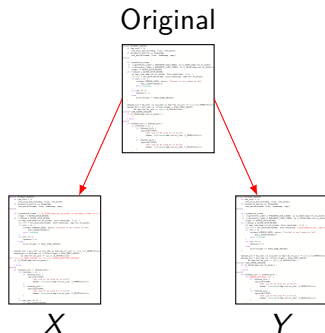
Main Result

Considering the scenario where

- File Y is an edited version of file X ,
- Alice has file X ,
- Bob has file Y ,
- Alice $\xleftrightarrow{\text{Comm. Channel}}$ Bob,

we provide a synchronization scheme with

- little communication required (within a constant factor of optimal rate),
- negligible error probability.



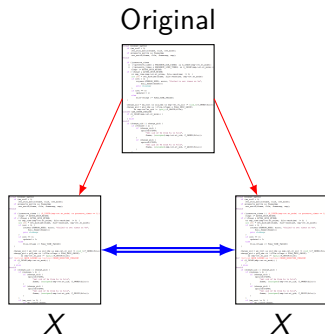
Main Result

Considering the scenario where

- File Y is an edited version of file X ,
- Alice has file X ,
- Bob has file Y ,
- Alice $\xleftrightarrow{\text{Comm. Channel}}$ Bob,

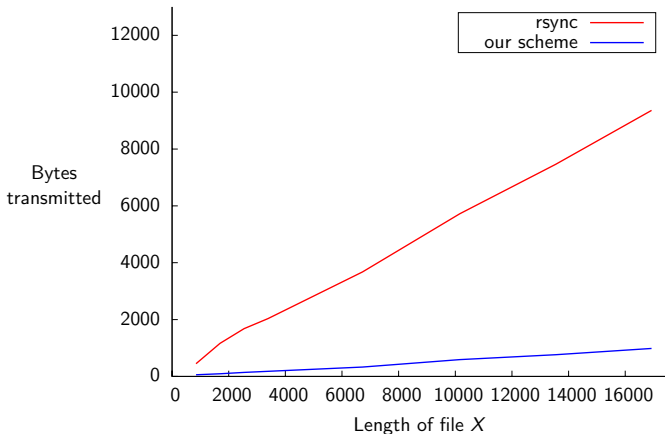
we provide a synchronization scheme with

- little communication required (within a constant factor of optimal rate),
- negligible error probability.



Comparison with rsync

- $\beta = 0.01$, i.i.d. file, i.i.d. edits,
- rsync block-size optimized,
- Our pivot length: 5, segment length: 100.



- Allow for more complex edit patterns (e.g., in practical scenarios, edits are often “bursty”),

- Allow for more complex edit patterns (e.g., in practical scenarios, edits are often “bursty”),
- Specialize our scheme to application-dependent types of files (e.g., if the files are source code, exploit that structure),

- Allow for more complex edit patterns (e.g., in practical scenarios, edits are often “bursty”),
- Specialize our scheme to application-dependent types of files (e.g., if the files are source code, exploit that structure),
- Optimize the implementation (in terms of both computation and network usage).