



Modular Multimodal Machine Learning for Extraction of Theorems and Proofs in Long Scientific Documents

Shrey Mishra, [Antoine Gauquier](#), [Pierre Senellart](#)

Motivation for TheoremKB

Volume of mathematical papers

Number of Arxiv papers per month since 1991

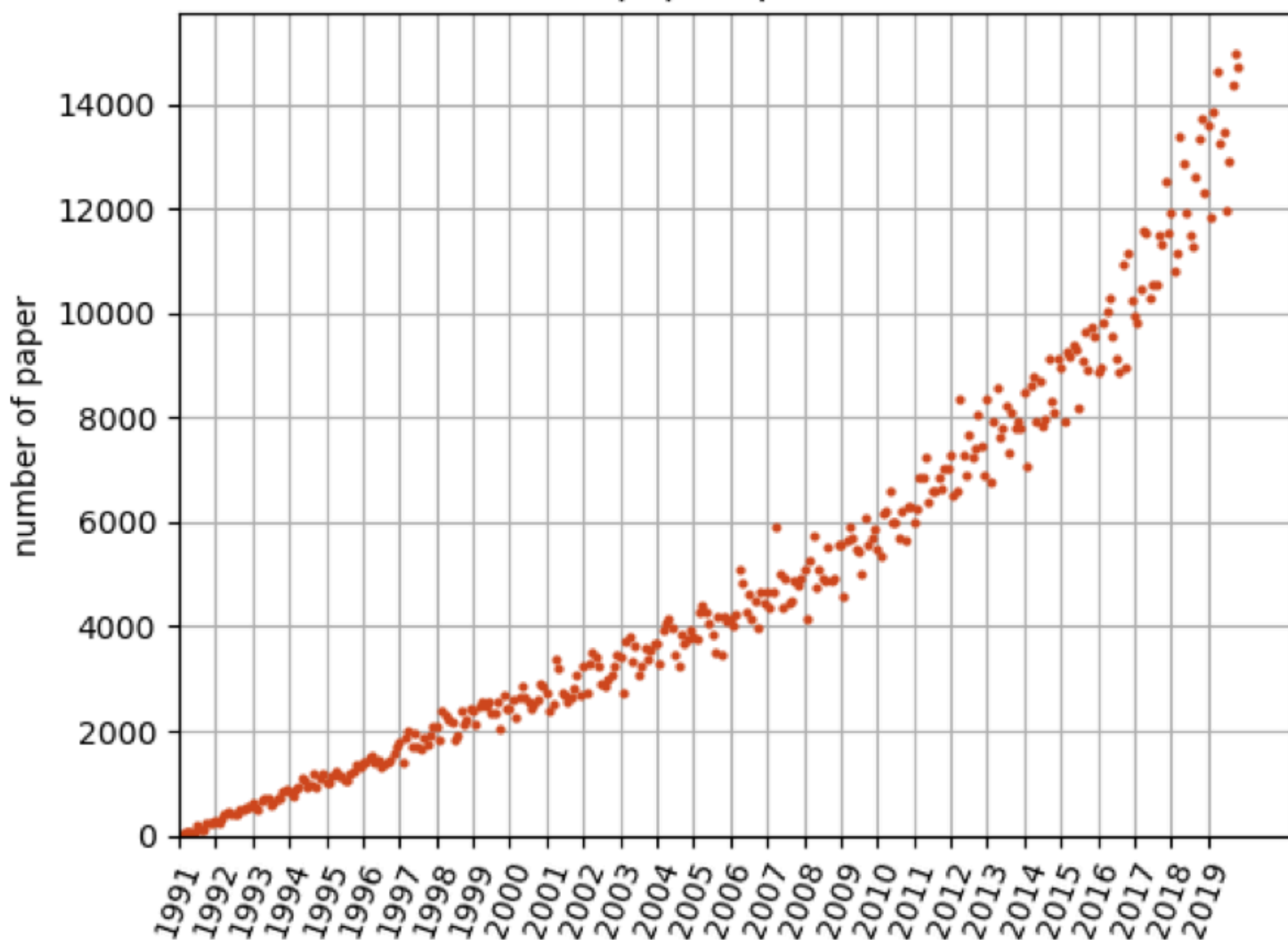


Figure 3.1: Total number of papers published on arXiv until 2019 [Del20]

Proportion of selected papers per month

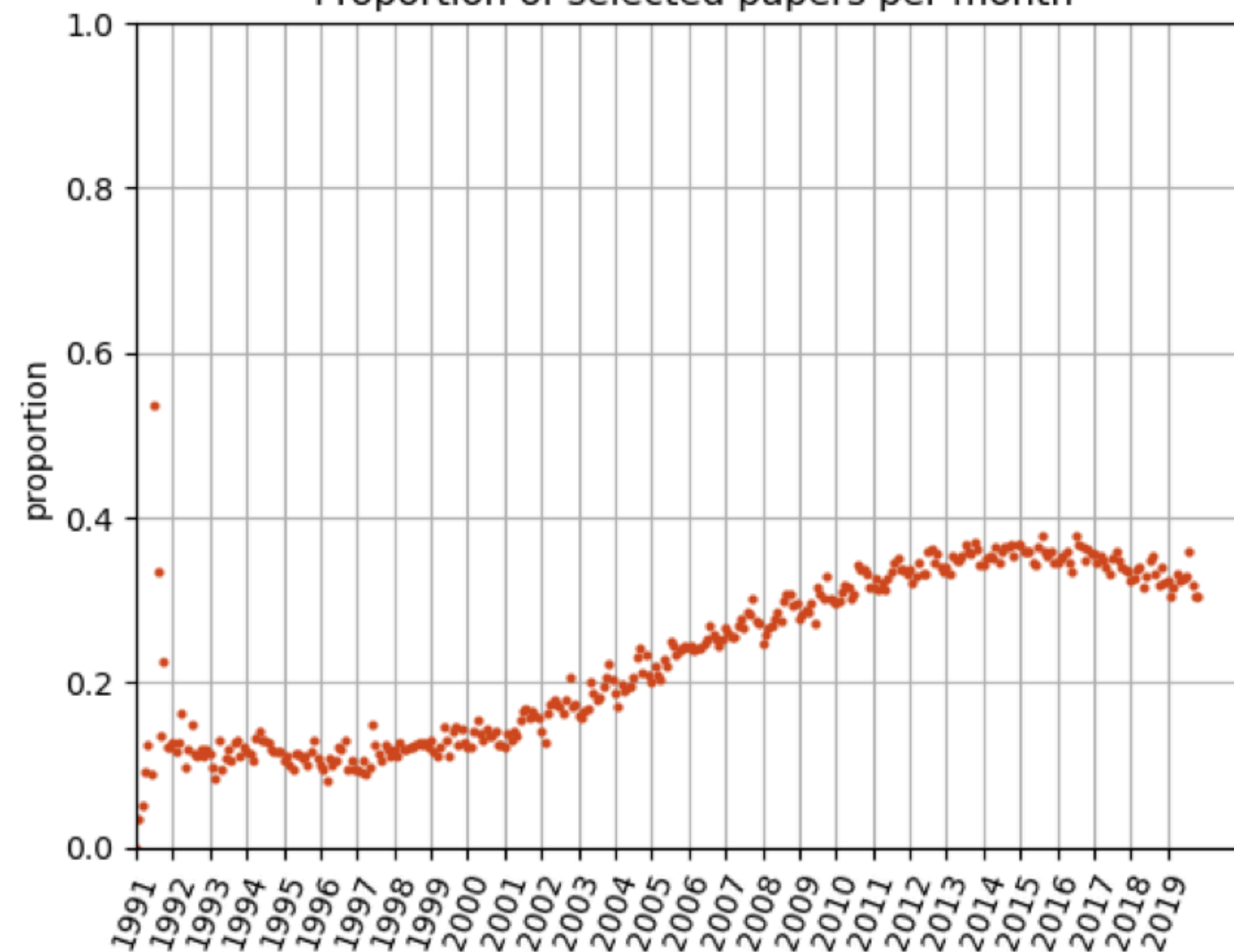


Figure 3.2: Papers with mathematical information such as Theorem, a Lemma or a Proposition [Del20]

Discovering existing knowledge



Michael et al. 2019

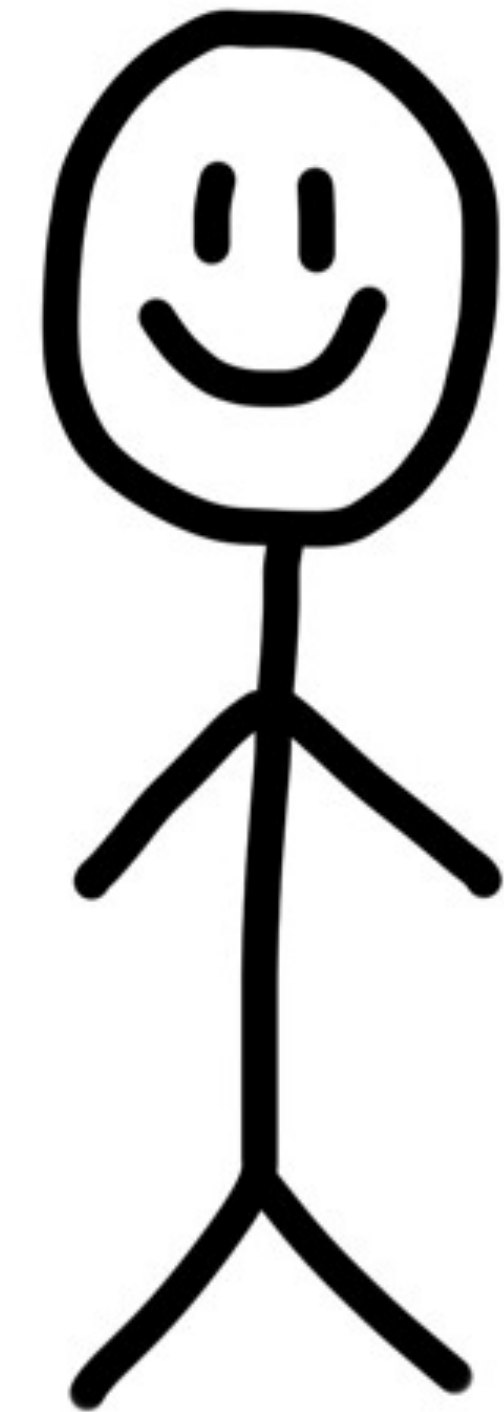


Hernich et al. 2012

Proves a result on Termination for linear TGDs

This other paper existed before

Michael



(likes to read mathematical papers)

Correcting publicised errors

Proof: Dichotomy for evaluating conjunctive queries



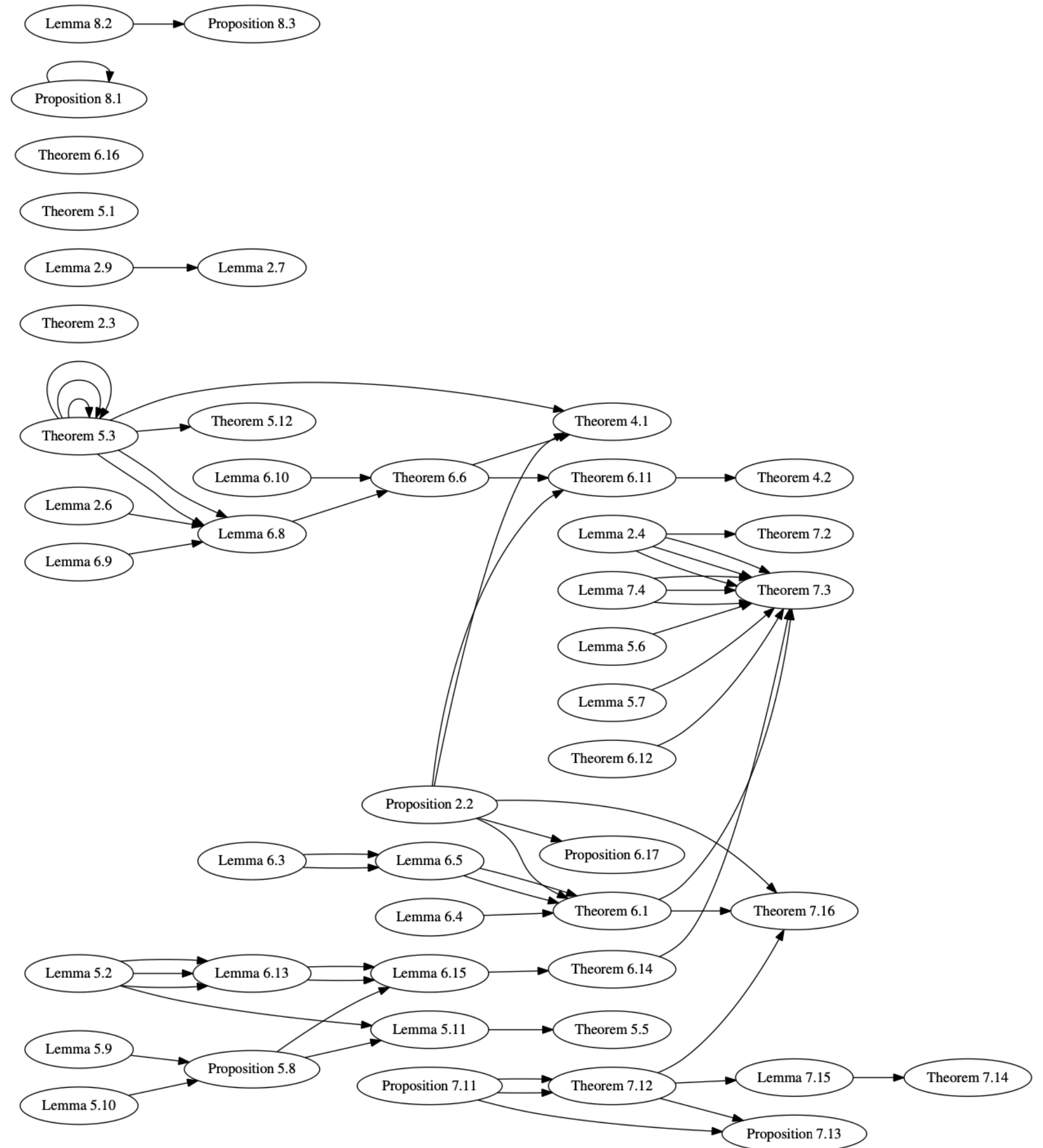
Nilesh et al. 2006

Proof v2: Dichotomy for evaluating conjunctive queries



Nilesh et al. 2012

Managing Complex Dependencies



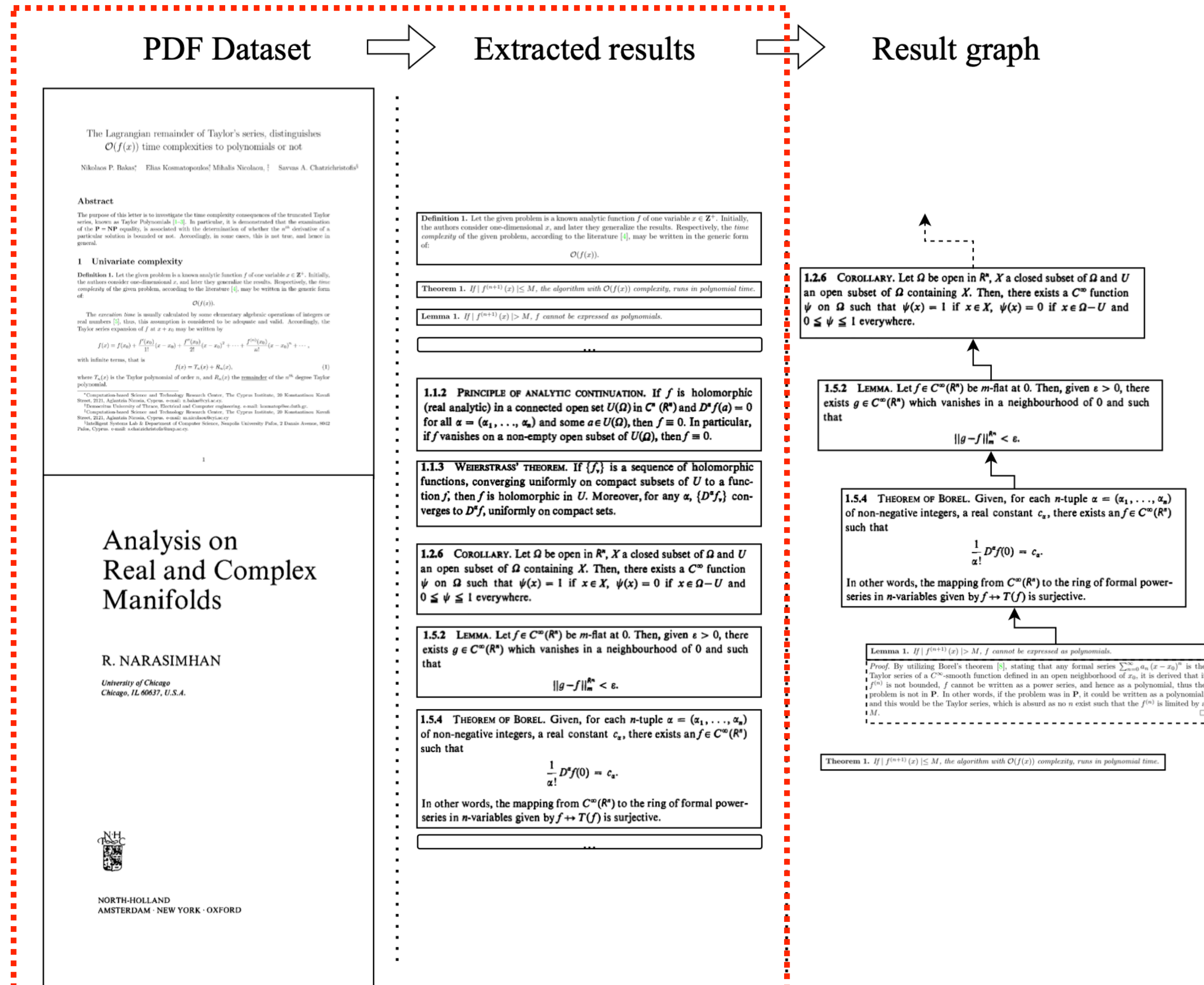
Senellart et al. Theory of Computing systems, 2019

Table of contents

- *Motivation of theoremkb*
- Information Extraction task
- Unimodal and Multimodal backbones
- Sequential approach

Extraction task

High level overview of TheoremKB



What is Extraction?

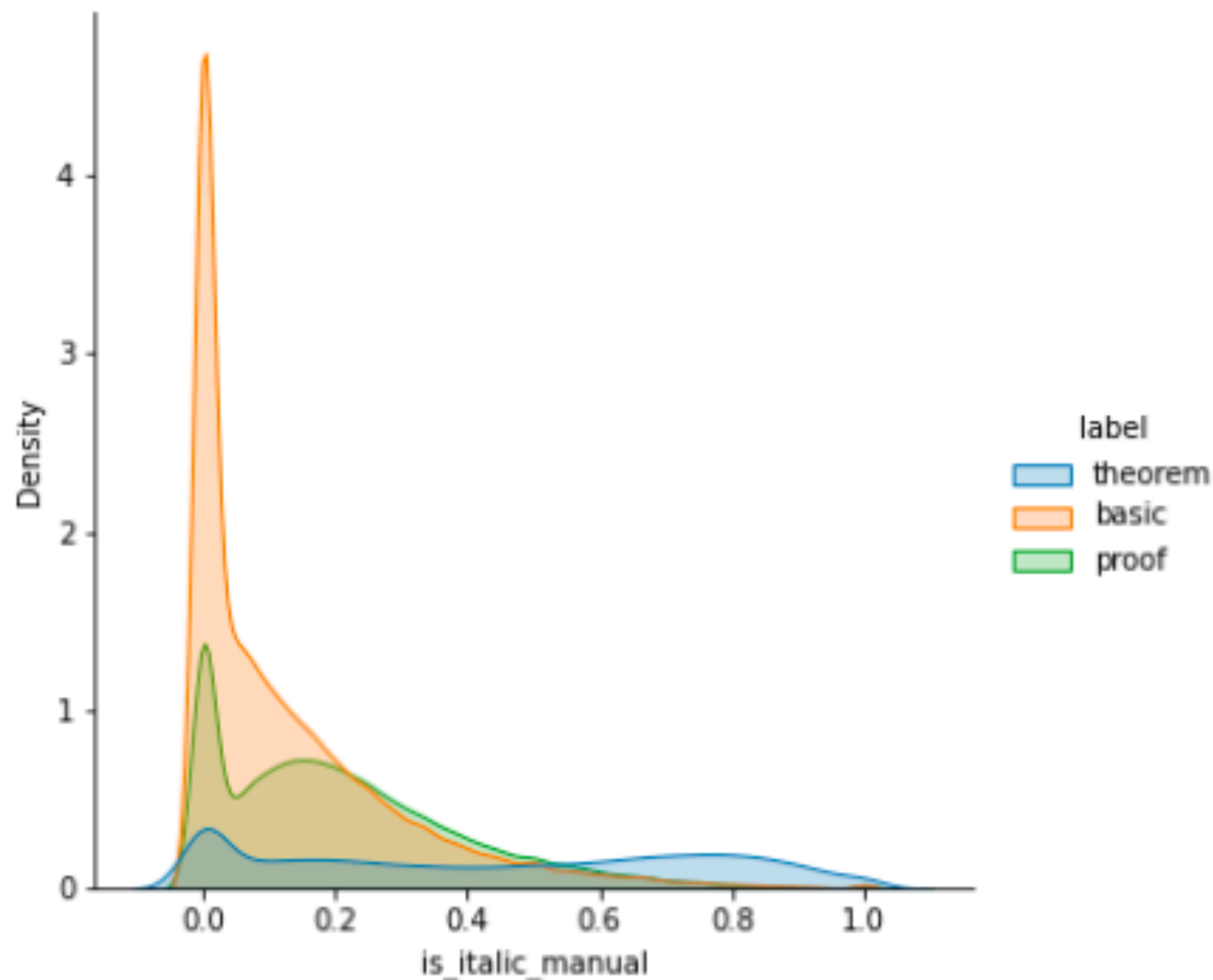
Theorem 1.8 ([HRVW09]). Let $f : \{0, 1\}^n \rightarrow \{0, 1\}^n$ be a one-way function, let X be uniformly distributed in $\{0, 1\}^n$, and let (Y_1, \dots, Y_m) be a partition of $Y = f(X)$ into blocks of length $O(\log n)$. Then (Y_1, \dots, Y_m, X) has next-block accessible entropy at most $n - \omega(\log n)$.

Proof. Since f is (t, ε) -one-way, the distributional search problem $(\Pi^f, f(X))$ where $\Pi^f = \{(f(x), x) : x \in \{0, 1\}^n\}$ is (t, ε) -hard. Clearly, $(f(X), X)$ is supported on Π^f , so by applying Theorem 3.8, we have that $(\Pi^f, f(X), X)$ has witness hardness $(\Omega(t), \log(1/\varepsilon))$ in relative entropy and $(\Omega(t), \log(1/\varepsilon) - \log(2/\delta))$ in $\delta/2$ -min relative entropy. Thus, by Theorem 4.7 we have that $(Y_1, \dots, Y_{n/\ell}, X)$ has next-block inaccessible relative entropy $(\Omega(t \cdot \Delta \cdot \ell^2 / (n^2 \cdot 2^\ell)), \log(1/\varepsilon) - \Delta)$ and next-block inaccessible δ -min relative entropy $(\Omega(t \cdot \delta \cdot \Delta \cdot \ell^2 / (n^2 \cdot 2^\ell)), \log(1/\varepsilon) - \log(2/\delta) - \Delta)$, and we conclude by Theorem 4.9.

□

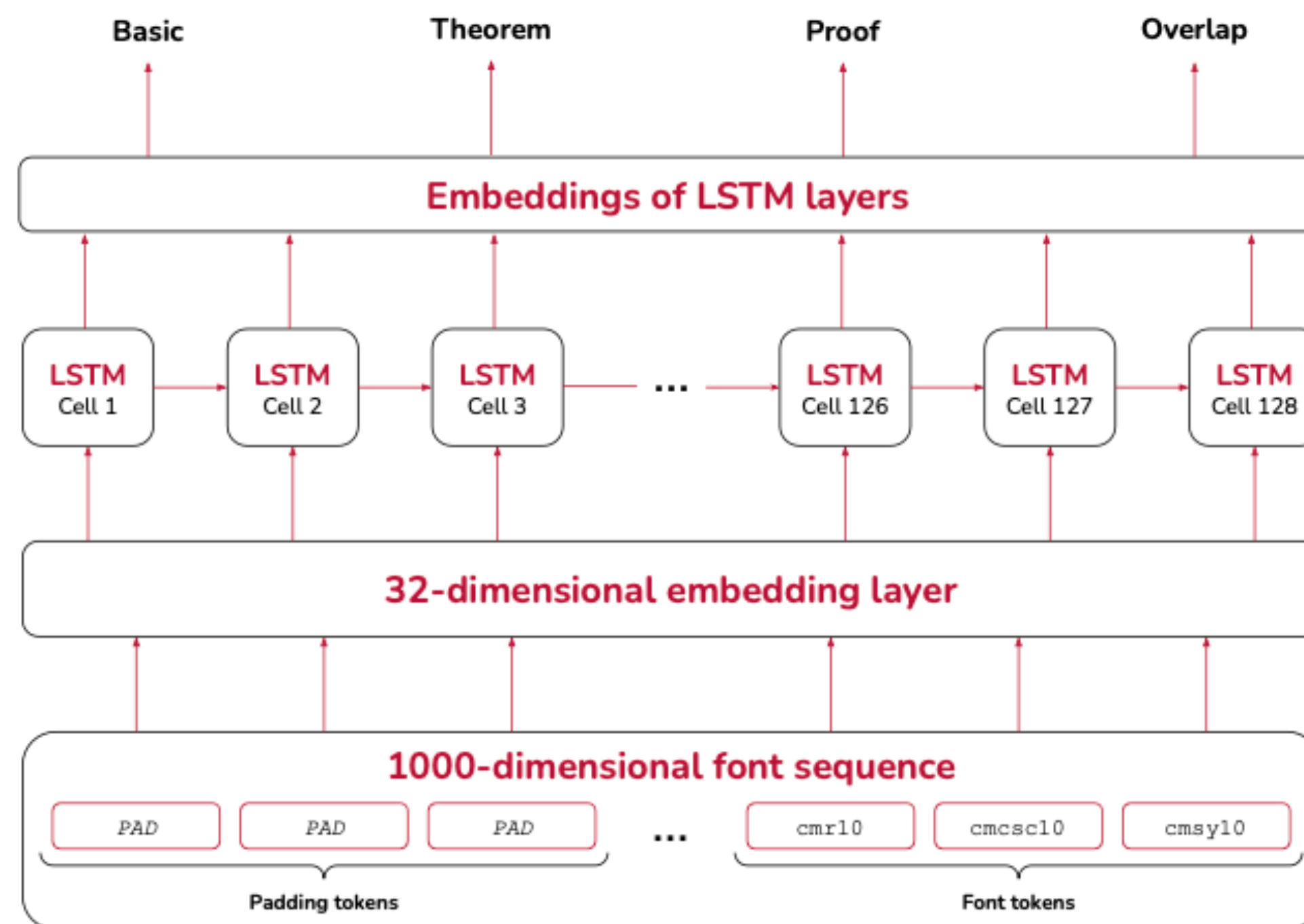
Unimodal backbones

Font based



Theorem 1.8 ([HRVW09]). Let $f : \{0, 1\}^n \rightarrow \{0, 1\}^n$ be a one-way function, let X be uniformly distributed in $\{0, 1\}^n$, and let (Y_1, \dots, Y_m) be a partition of $Y = f(X)$ into blocks of length $O(\log n)$. Then (Y_1, \dots, Y_m, X) has next-block accessible entropy at most $n - \omega(\log n)$.

Proof. Since f is (t, ε) -one-way, the distributional search problem $(\Pi^f, f(X))$ where $\Pi^f = \{(f(x), x) : x \in \{0, 1\}^n\}$ is (t, ε) -hard. Clearly, $(f(X), X)$ is supported on Π^f , so by applying Theorem 3.8, we have that $(\Pi^f, f(X), X)$ has witness hardness $(\Omega(t), \log(1/\varepsilon))$ in relative entropy and $(\Omega(t), \log(1/\varepsilon) - \log(2/\delta))$ in $\delta/2$ -min relative entropy. Thus, by Theorem 4.7 we have that $(Y_1, \dots, Y_{n/\ell}, X)$ has next-block inaccessible relative entropy $(\Omega(t \cdot \Delta \cdot \ell^2 / (n^2 \cdot 2^\ell)), \log(1/\varepsilon) - \Delta)$ and next-block inaccessible δ -min relative entropy $(\Omega(t \cdot \delta \cdot \Delta \cdot \ell^2 / (n^2 \cdot 2^\ell)), \log(1/\varepsilon) - \log(2/\delta) - \Delta)$, and we conclude by Theorem 4.9. \square

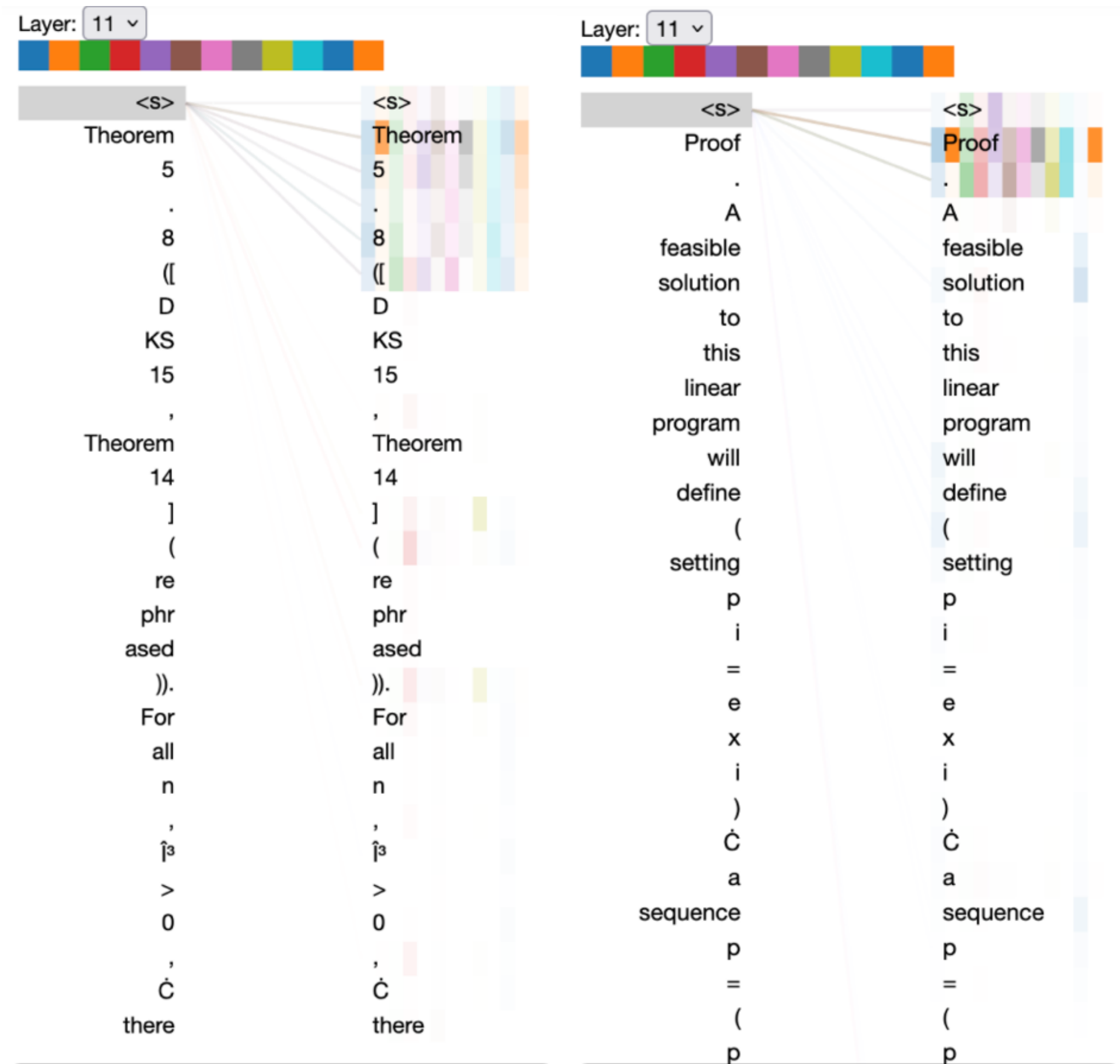


Vision Based

Proof. By Theorem 3, relative to a suitable oracle $A_{\mathcal{D}}$ (in fact, a *random* oracle suffices), there exists a signature scheme \mathcal{D} , such that any quantum chosen-message attack against \mathcal{D} must make superpolynomially many queries to $A_{\mathcal{D}}$. The oracle $A_{\mathcal{S}}$ will simply be a concatenation of $A_{\mathcal{M}}$ with $A_{\mathcal{D}}$. Relative to $A_{\mathcal{S}}$, we claim that the mini-scheme \mathcal{M} and signature scheme \mathcal{D} are *both* secure—and therefore, by Theorem 16, we can construct a secure public-key quantum money scheme \mathcal{S} .

Theorem 4. *Suppose the ETH holds with constant c . Then for every $\alpha, \beta \in \mathbb{N}$ there exists a $\gamma = O(\alpha + \beta)$ such that*

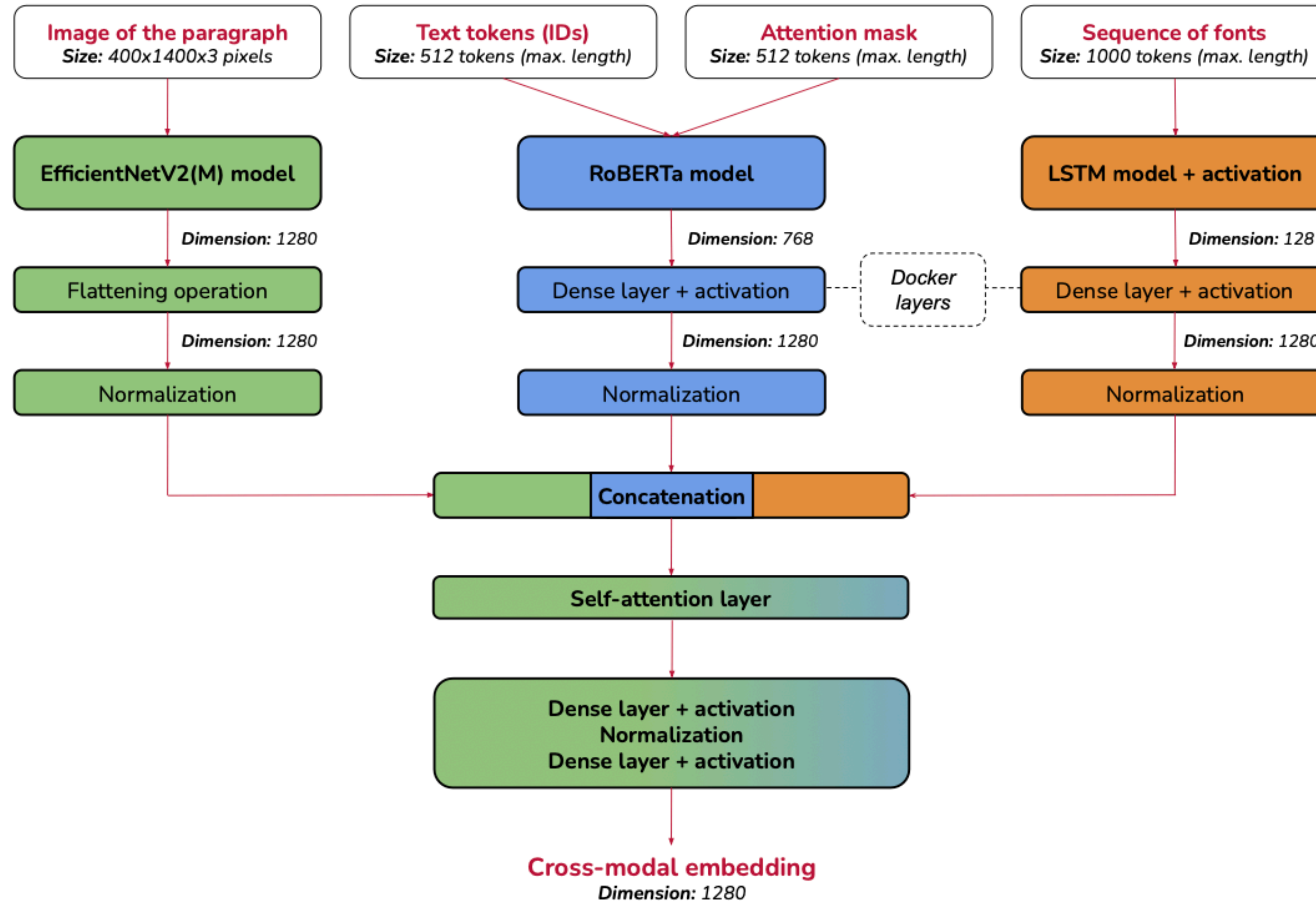
Text based



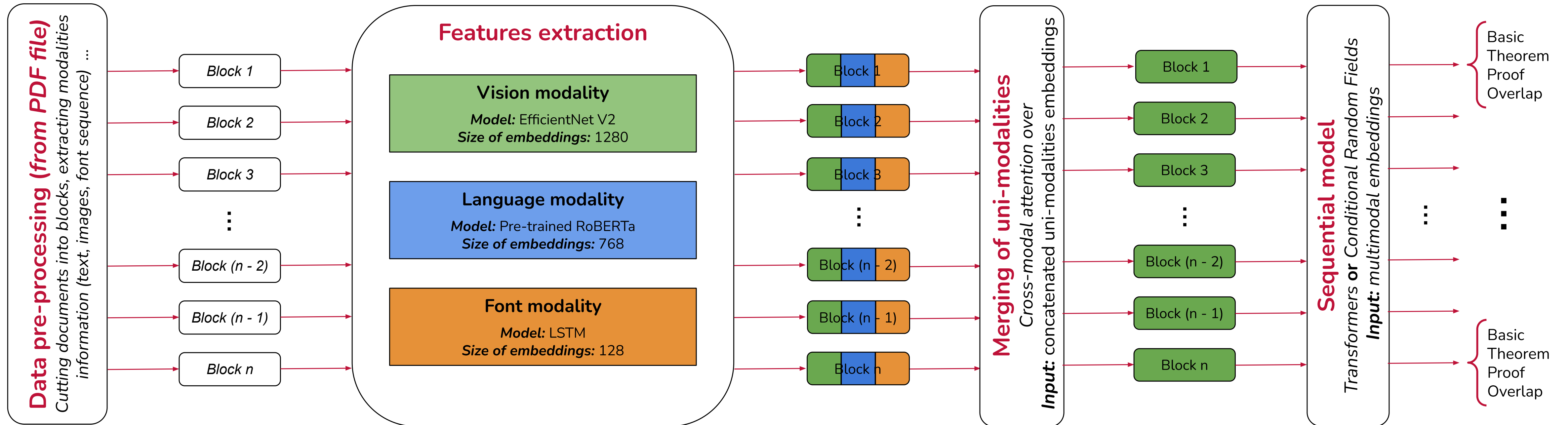
Theorem 5.8 ([DKS15, Theorem 14] (rephrased)). For all $n, \gamma > 0$, there exists a set $S \gamma \subseteq \text{PBD } n$ such that:

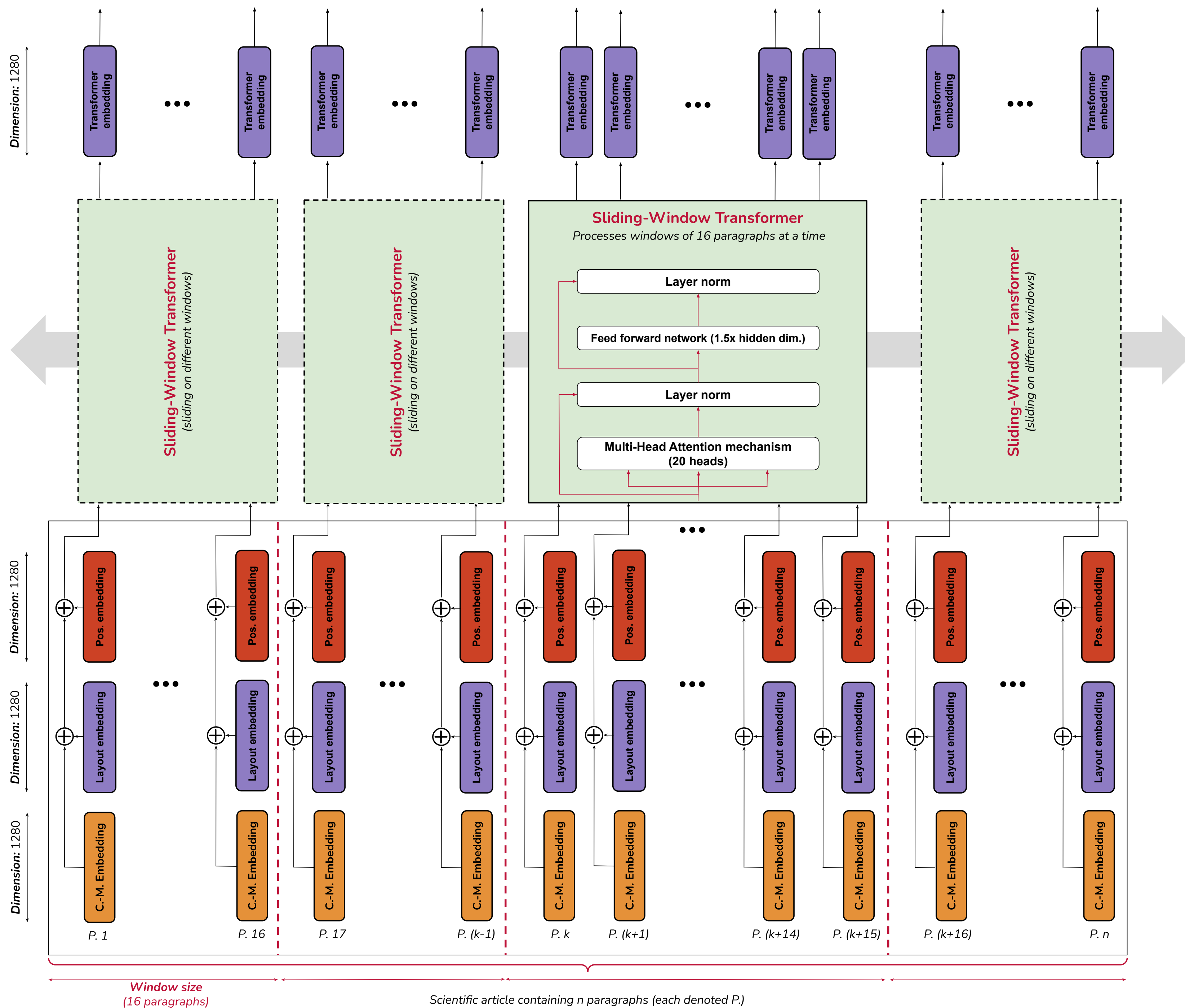
Proof. A feasible solution to this linear program will define (setting $p_i = e x_i$) a sequence $p = (p_1, \dots, p_n) \in (0, 1]^n$ such that

Multimodal model (raw features)



Sequential Approach





Performance of Sequence model

Modality	Model chosen	Seq. approach	#Batches	#Params (M)	Accuracy (%)	Mean F ₁ (%)
Dummy	always predicts <i>basic</i>	—	—	—	59.41	24.85
Top- <i>k</i> first word	use only first word	—	—	—	52.84	44.20
Line-based [MPS21]	Bert (fine-tuned)	—	—	110	57.31	55.71

Closing remarks

- Our multimodal approach can be adapted to long documents
- Can make inference on entire pdf in a single forward pass
- Comparable, consistent and computationally efficient
- Unlike many other approaches that rely on an OCR preprocessing to be useful (LayoutLM) ours rely on Grobid which is many times faster
- Our approach is Fast (encoder only) and Scalable and applicable in real world (~200k pdfs tested)
- Our approach captures cross modality without adding special losses