

LLM-PF: A LLM-Based Framework for Cryptographic Protocol Flow Extraction via a Novel IR towards Formal Verification

Anonymous ACL submission

Abstract

Formal verification is a robust and proven method for ensuring cryptographic protocol security, yet constructing formal models from natural language specifications remains a labor-intensive task demanding specialized expertise. While Large Language Models excel at text understanding, they struggle to directly generate accurate formal models due to the intricate logic and ambiguities in protocol descriptions. To bridge this gap, we propose a novel intermediate representation based on finite state machines. This representation deconstructs protocol logic into modular components, enabling LLMs to extract information step-by-step while ensuring a seamless and complete conversion to formal models. Building on this design, we developed a comprehensive extraction framework to automate the process. Extensive evaluations on a diverse set of protocols, ranging from basic to large-scale industrial standards, demonstrate the effectiveness and potential of our approach in automating formal verification.

1 Introduction

Formal analysis techniques enable rigorous security verification of cryptographic protocols that underpin secure communication systems by constructing precise models and applying mathematical reasoning (Goldreich, 2003; Meadows, 2003). Their reliability critically depends on accurately translating natural-language protocol specifications into unambiguous formal representations. However, such translation is inherently challenging, as protocol specifications are often written in informal prose, describe complex protocol flows and state transitions, and contain implicit assumptions that are difficult to formalize.

Given the difficulty of this translation task, recent advances in Large Language Models (LLMs) have attracted increasing attention as a potential means of assisting protocol understanding and formal modeling. LLMs have demonstrated strong

capabilities in text comprehension, and recent studies show that they can help interpret cryptographic protocol specifications and generate candidate formal models for subsequent analysis (Mao et al., 2025; Curaba et al., 2024; Duclos et al., 2024). Despite this promise, existing LLM-based approaches remain limited in practice.

One key limitation arises from the structural characteristics of large-scale cryptographic protocol standards. Such documents are lengthy and informationally dispersed, with non-sequential dependencies that scatter protocol logic across multiple sections. Existing work attempts to cope with this complexity either by relying on simplified descriptions, such as Alice&Bob narratives (Basin et al., 2015; Keller and Basin, 2014; Curaba et al., 2024; Li et al., 2025), or by designing specialized pipelines for direct document analysis (Mao et al., 2025; Duclos et al., 2024). While effective for small protocols, both strategies struggle to scale to large and complex specifications.

A second limitation concerns factual hallucination stemming from LLMs' lack of domain-specific protocol knowledge (Huang et al., 2025). When processing large-scale protocol specifications with rich state spaces and intricate execution paths, LLMs may introduce incorrect assumptions, omit critical conditions, or generate inconsistent transitions. Although mitigation techniques such as fine-tuning or repair modules have been explored, hallucination remains a persistent issue, particularly in settings that require precise reasoning about protocol states and message dependencies (Mao et al., 2025).

To address these challenges, we propose LLM-PF, an LLM-driven framework for extracting protocol flows from natural-language specifications using an Intermediate Representation (IR). Central to our approach is the Protocol Finite State Machine (PFSM), a finite-state-machine-based representation that explicitly captures protocol roles,

084 states, transitions, and message dependencies. By
085 structuring protocol knowledge into the PFSM, our
086 framework decouples the extraction process from
087 downstream formal verification tasks. This design
088 significantly mitigates hallucination and enables
089 the scalable analysis of large, complex protocols.

090 Our approach draws inspiration from prior work
091 on extracting state machines from network proto-
092 cols (Sharma and Yegneswaran, 2023; Miao et al.,
093 2025; Pacheco et al., 2022). However, adapting
094 these methods to cryptographic protocols requires
095 addressing fundamental distinctions (Hazzazi et al.,
096 2023). Unlike network protocol state machines,
097 which primarily focus on state transitions, our
098 PFSM must explicitly describe cryptographic oper-
099 ations and their interaction relationships. This addi-
100 tional semantic complexity is essential for rigorous
101 formal analysis and distinguishes our representa-
102 tion from conventional network models.

103 The main contributions of this paper are summa-
104 rized as follows:

- 105 • **We introduce PFSM**, a novel intermediate
106 representation based on Finite State Machine
107 (FSM) that is LLM-friendly and supports loss-
108 less translation into formal verification lan-
109 guages such as ProVerif (Blanchet, 2012) and
110 Tamarin (Meier et al., 2013). We further re-
111 lease a curated dataset of 10 real-world cryp-
112 tographic protocols annotated in PFSM.
- 113 • **We develop an automated extraction frame-**
114 **work** that leverages PFSM to process large-
115 scale protocol specifications. The framework
116 integrates workflow decomposition and en-
117 hanced retrieval to improve extraction accu-
118 racy, with experimental results demonstrating
119 its effectiveness on large and complex proto-
120 cols.
- 121 • **We design a toolchain** that automatically
122 converts PFSM into formal verification code
123 and protocol sequence diagrams (e.g., Plan-
124 tUML (Roques and Contributors, 2025) and
125 Mermaid (Sveidqvist and Contributors to Mer-
126 maid, 2014)).

127 We open-source our dataset and code for future
128 research and reproducibility ¹.

¹[https://anonymous.4open.science/r/
Anonymous-LLMPF-975D](https://anonymous.4open.science/r/Anonymous-LLMPF-975D)

2 Related Works 129

We review related work from the perspective of 130
how protocol specifications are translated toward 131
formal analysis. In particular, prior research differs 132
in both the form of input it assumes and the tech- 133
niques it aims to achieve. Based on these dimen- 134
sions, existing approaches can be broadly grouped 135
into three categories, which we discuss in the fol- 136
lowing subsections. 137

Direct Conversion from Simplified Inputs. 138

This line of work focuses on translating high-level, 139
structured descriptions into formal models. Recent 140
LLM-based approaches like CryptoFormalEval and 141
P2FGPT convert Alice&Bob style inputs into mod- 142
els for Tamarin and ProVerif (Curaba et al., 2024; 143
Li et al., 2025). Earlier work, such as Keller’s 144
A&B language, follows a similar structured input 145
paradigm (Basin et al., 2015; Keller and Basin, 146
2014), while approaches using simple natural lan- 147
guage descriptions have also been explored (Duclos 148
et al., 2024). However, these methods suffer from 149
poor scalability and generalization when facing the 150
complexity and ambiguity of real-world protocol 151
RFCs. 152

IRs for Protocol Analysis and Verification. 153

This line of work addresses protocol complexity by 154
extracting intermediate representations from RFCs. 155
A common approach is to derive structured mod- 156
els, such as finite state machines (FSMs), for tasks 157
including vulnerability analysis and attack synthe- 158
sis (Sharma and Yegneswaran, 2023; Goo et al., 159
2019; Al Ishtiaq et al., 2024). 160

These IRs are typically designed to capture high- 161
level state transitions while abstracting away cryp- 162
tographic operations and message-level seman- 163
tics (Hazzazi et al., 2023). As such, they are well- 164
suited for exploratory analysis, but are not intended 165
to directly support automated formal verification. 166
Nevertheless, the role of IRs as a bridge toward 167
formal reasoning is widely recognized. For exam- 168
ple, MetaCP introduces an XML-based IR that can 169
be accurately translated into multiple formal tools, 170
although it relies on user-provided designs rather 171
than parsing RFCs (Arnaboldi and Metere, 2019). 172

End-to-End RFC to Formal Model via IR. 173

This line of work focuses on constructing end-to- 174
end pipelines that translate RFCs into formal mod- 175
els via an intermediate representation. For example, 176
Mao et al. (2025) proposed an LLM-assisted ap- 177
proach based on a lambda-calculus IR or translation 178
into SAPIC+. 179

3 PFSM: An FSM-Based Protocol IR

To bridge the gap between the linguistic capabilities of LLMs and the precision required by formal models, we introduce PFSM, a FSM-based IR that balances comprehensibility for LLM-assisted extraction with convertibility to formal models.

3.1 Design Principles

The architecture of PFSM is guided by three core tenets, which collectively address the limitations of prior work:

Expressive Power through State-Driven Modeling: PFSM harnesses the intrinsic state-transition mechanism of FSMs to natively represent the intricate logic of cryptographic protocols, including conditional branching and complex flows, thus scaling to sophisticated systems.

Scalable Extraction via Modular Decomposition: We partition protocols into atomic state units, transforming monolithic descriptions into a modular architecture. This enables a divide-and-conquer strategy for LLMs to process states incrementally, mitigating context-window limits and enhancing fidelity for large-scale protocols.

Seamless Convertibility through Structural Congruence: PFSM semantics align with formal verifiers like ProVerif and Tamarin. States encapsulate knowledge and actions, while transitions dictate evolution. This ensures a direct, structurally isomorphic mapping that guarantees conversion completeness and soundness.

3.2 Core Syntax and Design Rationale

At its core, a protocol PFSM \mathcal{P} is defined as:

$$\mathcal{P} = (\mathcal{E}, Ks, M) \quad (1)$$

where \mathcal{E} is the finite set of participating entities, Ks specifies the initial knowledge distribution, and M assigns a finite state machine to each entity.

The initial knowledge distribution function:

$$Ks : \mathcal{K}_{init} \mapsto 2^{\mathcal{E}} \quad (2)$$

maps each piece of initial knowledge to the subset of entities that process it. This explicit representation provides a clear blueprint for protocol initialization and supports consistent translation into formal verification models. In practice, such knowledge assignments can be instantiated as parameters, fresh values, or persistent facts in mainstream formal tools, reducing ambiguity introduced by informal specifications.

Each entity $e \in \mathcal{E}$ is associated with a unique finite state machine FSM_e , defined by the assignment function $M : \mathcal{E} \mapsto FSM$. The behavior of each entity e is captured by:

$$FSM_e = \{FSM_{e,s} \mid s \in S_e\} \quad (3)$$

where S_e denotes the finite set of state identifiers for e .

Each state unit $FSM_{e,s}$ is a tuple:

$$FSM_{e,s} = (K_{e,s}, O_{e,s}, \theta_{e,s}) \quad (4)$$

where $K_{e,s}$ is the knowledge set held by entity e while in state s . $O_{e,s}$ denotes the set of actions executable upon entering s , and $\theta_{e,s}$ defines the transition rules to subsequent states. State transitions are triggered by message actions, or by decision-based branch conditions.

This action-driven granularity is a deliberate design choice. By organizing protocol behavior around explicit actions rather than informal narratives, PFSM enables modular decomposition of complex protocols into atomic, analyzable units. Such decomposition supports a divide-and-conquer extraction strategy, allowing LLMs to process protocol behavior incrementally and mitigating context-window limitations when handling large specifications.

Moreover, anchoring state transitions to explicit actions aligns PFSM with the operational abstractions used by formal verification tools. Sequences of actions within a state can be mapped to process interactions in the π -calculus for ProVerif (Abadi et al., 2017), while the encapsulated structure of state units aligns with the rule-based modeling style of Tamarin (Meier et al., 2013). This structural congruence supports systematic translation to formal code and reduces reliance on ad hoc interpretation during model construction.

Overall, the core syntax of PFSM balances expressiveness, modularity, and convertibility. By making protocol control flow, knowledge evolution, and actions explicitly at the state level, PFSM provides a representation that is well suited for LLM-Driven extraction and downstream formal verification.

Appendix A details the translation semantics from PFSM syntax to ProVerif, and provides a simple illustrative example.

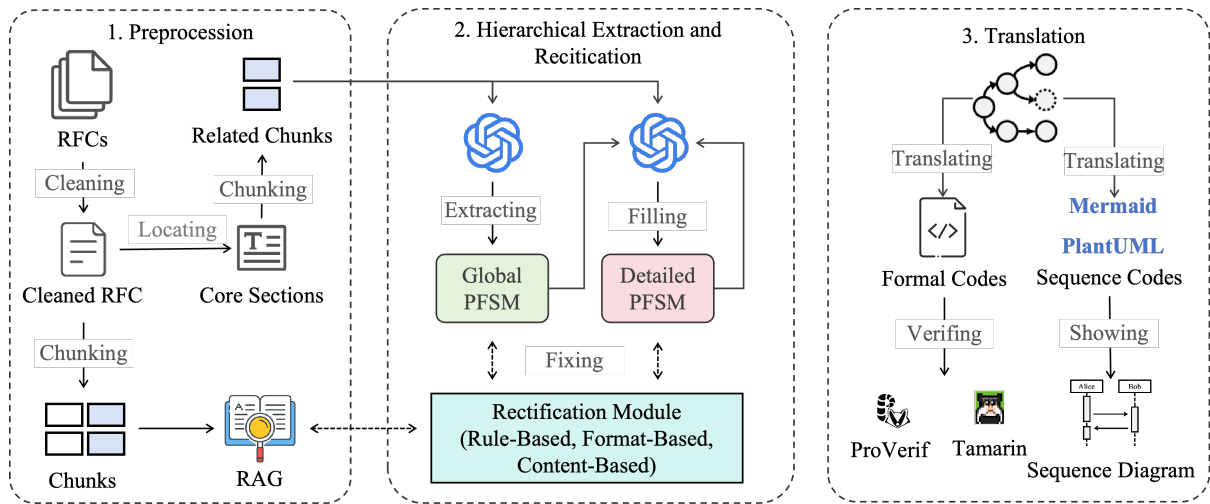


Figure 1: Overall framework of the LLM-PF

4 PFSM Extraction Framework

This chapter presents a systematic hierarchical extraction framework designed to automatically construct PFSMs from natural language specifications. The overall pipeline is illustrated in Figure 1:

4.1 Preprocessing

To address the challenges posed by lengthy and fragmented protocol specifications, which exceed the effective context window and comprehension capabilities of LLMs, a preprocessing phase is first employed to distill and structure the raw documents.

Standardization and Data Cleaning. We first employ document parsing tools to convert multi-source heterogeneous documents into a unified Markdown format. Subsequently, we utilize regular expressions combined with heuristic rules to eliminate explicit noise, such as page numbers, hyperlinks, and images, ensuring that the input data is clean and consistent.

Core Section Localization. Cryptographic protocol specifications typically exhibit a rigid structure where descriptions of protocol logic are concentrated in specific sections—for instance, the *Protocol Details* in UAF (FIDO Alliance, 2017) or the *Authentication Framework* in 5G (3GPP, 2022). Conversely, auxiliary sections such as the introduction, definitions, security proofs, and appendices, while providing context, are largely irrelevant to the control flow. Due to the limitations of LLMs in precisely locating information within massive unstructured text (Hengle et al., 2025; Liu et al., 2025b), our framework employs semantic recog-

nition to identify these core sections, designating them as primary input for PFSM generation.

Semantic Chunking. To address LLM context limits, we chunk localized protocol sections. Simultaneously, we construct a vector knowledge base via fine-grained chunking of the entire specification to mitigate cross-chapter information loss. This allows the subsequent rectification algorithm to retrieve missing contextual definitions via Retrieval-Augmented Generation (RAG) (Lewis et al., 2020), ensuring complete state descriptions.

4.2 Hierarchical Extraction and Rectification

Extracting PFSM from unstructured protocol specifications presents three key challenges for LLMs.

C1: Noise interference in entity identification. Protocols often contain non-functional components (e.g., PKI or databases) that models may misclassify as participants. Furthermore, physically isolated components described in specifications should be abstracted into a single logical entity, increasing model comprehension complexity. For instance, this abstraction is required in the 5G AKA protocol, where the internal USIM and ME components of a User Equipment (UE) should be merged into one entity (3GPP, 2022).

C2: Complexity of state transitions. Protocol operations involve multi-branch logic and fine-grained cryptographic primitives, making it difficult for models to accurately parse all atomic operations and state flows in a single pass.

C3: Difficulty in capturing inter-entity dependencies. The core of a PFSM lies in the inter-entity dependencies established via message exchange; in multi-branch scenarios, correctly aligning states

across different entities often exceeds the foundational capabilities of current models.

To address these challenges, we propose a hierarchical extraction framework that decomposes the task to maximize the model’s reasoning efficacy. The algorithm consists of four stages, as outlined in Algorithm 1.

Entity Identification (Addressing C1). This stage aims to filter non-functional noise and extract only participants \mathcal{E} directly involved in protocol interactions. For ambiguously defined entities, we additionally introduce a human-in-the-loop mechanism to calibrate the entity list via expert feedback, ensuring the accuracy of the basic modeling units.

Global PFSM Extraction (Addressing C2). Upon confirming the entity set, we construct a global PFSM \mathcal{P}_g . \mathcal{P}_g defines the macro-level state machine, specifying the name and description of each state for every entity $e \in \mathcal{E}$. It captures the state transition logic (i.e., the next state for a given entity) and explicitly encodes the inter-entity dependencies established via message exchange, specifically mapping each sending state to its corresponding receiving state. \mathcal{P}_g serves as a blueprint for subsequent refinement, clearly delineating the protocol’s topological structure.

Global PFSM Rectification (Addressing C3). To mitigate logical discontinuities in \mathcal{P}_g , we formulate a set of rules targeting specific integrity failures: *Content Integrity*, *Transition Integrity*, and *Interaction Consistency*. These rules explicitly identify violations—such as incomplete state fields, missing intermediate states, or unmatched interactions—and provide feedback to the LLM for iterative regeneration.

Detailed PFSM Filling and Rectification (Addressing C2). Based on the global skeleton \mathcal{P}_g , the algorithm proceeds to generate the detailed PFSM \mathcal{P} . To handle protocols that may exceed output token limits, we employ an iterative batch extraction strategy initialized with a batch size threshold t . As shown in Algorithm 1, the loop iteratively selects the next batch of states B via `KGetBatch`. Within each iteration, the `KFill` function constructs the detailed state information for B . This is immediately followed by the `KVerify` function, which initiates a verification and repair mechanism: rule-based checks validate the legality of atomic operations, and for any missing definitions or ambiguous parameters, RAG technology retrieves relevant information from the full document C to automatically complete the descriptions. The verified batch is

then aggregated into \mathcal{P}_d . This cycle repeats until all states are processed, culminating in the final, accurate PFSM.

Algorithm 1: Hierarchical Extraction and Rectification Algorithm

Input: Chunks C , Threshold t
Output: Verified PFSM \mathcal{P}

```

1  $\mathcal{E} \leftarrow \text{ExtractEntities}(C)$ 
2  $\mathcal{P}_g \leftarrow \text{GenGlobalPFSM}(C, \mathcal{E})$ 
3  $\mathcal{P}_g \leftarrow \text{RectifyStruct}(\mathcal{P}_g, C, \mathcal{E})$ 
4  $\mathcal{P}_d \leftarrow \emptyset$ 
5 while ( $B \leftarrow \text{NextBatch}(\mathcal{P}_g) \neq \emptyset$ ) do
6    $\mathcal{P}_t \leftarrow \text{FillDetails}(\mathcal{P}_g, B, C)$ 
7    $\mathcal{P}_t \leftarrow \text{VerifyAndFix}(\mathcal{P}_t, C)$ 
8    $\mathcal{P} \leftarrow \mathcal{P} \cup \mathcal{P}_d$ 
9 end
10  $\mathcal{P} \cup \mathcal{P}_d$ 
11 return  $\mathcal{P}$ 

```

4.3 Translation Module

To maximize the utility of the extracted PFSM, we developed a translation ecosystem supporting both formal verification and visualization. For rigorous security analysis, a high-performance translator converts the PFSM into specification languages (e.g., ProVerif) to verify properties like reachability (see Appendix A). Additionally, visualization translators (PlantUML, Mermaid) generate intuitive diagrams for inspection.

5 Experiments and Evaluation

This chapter presents a comprehensive evaluation of our proposed LLM-PF framework, guided by the following research questions:

- **RQ1:** How effectively does the LLM-PF framework extract PFSMs across protocols of varying scales?
- **RQ2:** To what extent does the LLM-PF framework outperform direct LLM baselines?
- **RQ3:** How do the hierarchical extraction and rectification mechanisms contribute to the validity and completeness of the generated models?

5.1 Dataset and Ground Truth

To ensure a robust analysis, we curated a dataset comprising 10 protocols ranging from

simple to complex benchmarks. The set centers on Needham-Schroeder Symmetric Key (NSSK) (1978) and Needham-Schroeder Public Key (NSPK) (1978), alongside their variants: Neuman-Stubblebine (Neuman) (1993), Yahalom (1990), Otway-Rees (OR) (1987), and Wide-Mouth-Frog (WMF) (1990). Beyond this group, we included the classic Kao-Chow (KC) (1995) protocol and three complex protocols: 5G AKA (2022), FIDO UAF (2017), and EDHOC (2021). Notably, for protocols lacking official RFC standards, we adapted specifications from the LLM-aided dataset (2025) and transformed them into RFC-style documents.

The ground-truth PFSM were manually constructed by domain experts based on original RFC documentation. For each protocol, we generated the corresponding formal verification code. To ensure consistency, these PFSMs were validated by automatically translating them into the written formal verification code. This rigorous methodology supports comprehensive end-to-end verification.

Table 1 provides an overview of the dataset, detailing the scale of the RFC documents, PFSMs, and the generated formal verification code.

Table 1: Overview of the Dataset Complexity

	Description		Rows		
	Roles	Rounds	RFCs	PFSM	ProVerif
NSSK (1978)	3	5	112	441	90
NSPK (1978)	2	3	88	296	64
NS (1993)	3	4	83	425	114
Yahalom (1990)	3	4	123	402	136
OR (1987)	3	4	99	433	118
WMF (1990)	3	2	66	252	89
KC (1995)	3	4	20	78	97
5G (2022)	4	11	470	1286	400
UAF (2017)	3	8	4010	1370	149
EDHOC (2021)	2	4	1428	3426	553

5.2 Evaluation Metrics

Given the extraction algorithm, we evaluate the LLM-PF in two dimensions.

Global State Machine Metrics. To evaluate the large model’s capability in recognizing cryptographic protocol state transitions and to validate the effectiveness of our global state extraction algorithm, we employ standard metrics for the extraction task. We count the correctly predicted states, the false positives predicted in error, and the ground truth states that were missed. Based on these counts, we report Precision (P_g), Recall

(R_g), and F1-score (F_g) as the primary evaluation metrics.

State Detail Metrics. To conduct a fine-grained assessment of the correctly identified states, we evaluate the fidelity of their internal attributes: the knowledge set K_s , the action set O_s , and the transition θ_s . For each attribute, we individually count the true matches, false positives, and false negatives. Specifically, knowledge correctness is assessed via the cosine similarity of textual descriptions followed by global substitution alignment; actions are evaluated by strict matching against defined function names and parameters; and transitions are determined by strict matching against defined transition rules. We calculate the Precision (P_d), Recall (R_d), and F1-score (F_d) for each attribute based on these counts and derive the overall metrics for the state details by taking the weighted average of these scores.

5.3 Results

We evaluated our method on four state-of-the-art LLMs: GPT-5.2 (OpenAI, 2025), Gemini-2.5 Pro (Comanici et al., 2025), DeepSeek-V3 (Liu et al., 2025a), and Grok-4 (xAI, 2025).

5.3.1 Overall Performance (RQ1)

Figure 2 illustrates the performance of the LLM-PF framework across a spectrum of cryptographic protocols, categorized by their complexity. For each protocol, the left subplot reports the P_g , R_g , and F_g for global state extraction, while the right subplot details the metrics P_d , R_d and F_d for detailed state extraction.

Validation and Generalization. The evaluation includes both lightweight protocols and industrial standards. For lightweight protocols, the extraction pipeline achieves precision consistently above 80%. While these results validate the fundamental soundness of the LLM-PF architecture on self-contained specifications, the primary objective of this work is to address the challenges posed by real-world industrial protocols. In this domain, the framework maintains moderate precision. This performance variation highlights the complexity of industrial-grade formalization, a task that remains challenging even for state-of-the-art models.

Findings on Industrial Protocol Complexity. Qualitative analysis reveals that performance on large protocols is constrained by both the inherent instability of LLMs and the structural ambiguity of technical specifications.

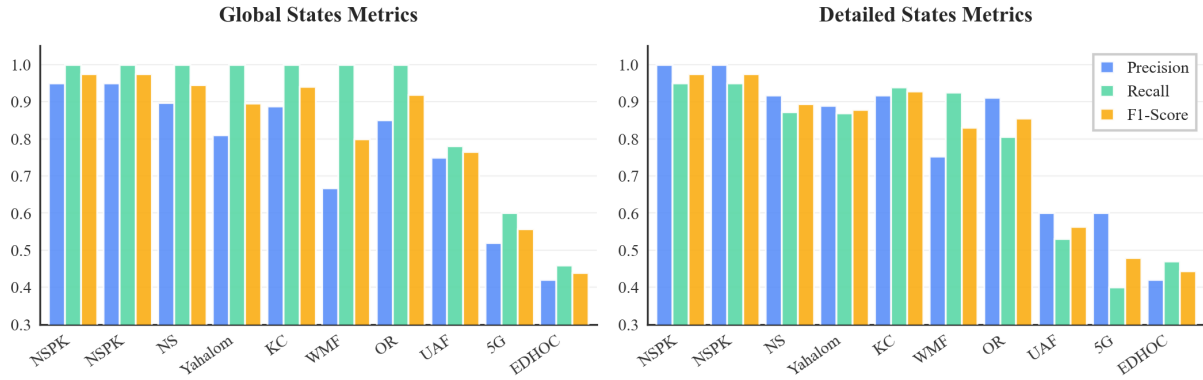


Figure 2: Overall Performance of LLM-PF across different protocols using the ChatGPT model.

Finding 1: Non-deterministic Generation. Although LLM-PF incorporates mechanisms to mitigate hallucinations, we observe notable variance in output quality across multiple runs, even with a temperature setting of 0.

Finding 2: Specification Structural Ambiguity. First, *cross-references* fragment logic; for instance, 5G AKA (p.52) defers AUTN verification to an external document (“as described in TS 33.102”), hindering the model’s ability to infer the complete mechanism (3GPP, 2022). Second, *noisy examples* mislead extraction; in UAF (Section 3.4.1), concrete registration request examples cause the model to erroneously treat instance-specific data as part of the generic flow (FIDO Alliance, 2017). Third, *long-range dependencies* make error handling difficult to capture. In 5G AKA, exception handling is isolated in a “Void” section (6.1.3.2.1, p.54) (3GPP, 2022). However, our results confirm that the Content Inspection and Repair Mechanisms are effective in bridging these dependencies.

Finding 3: Combinatorial Branching. A distinct challenge arises from branching logic, as seen in EDHOC (Selander et al., 2021), which defines multiple mutually exclusive authentication modes. This complexity leads to a precision of approximately 40%, as the model struggles to maintain a coherent global state across disjointed branches. This suggests that future work should explore protocol pre-processing strategies, such as extracting different modes separately.

Model Capability Analysis. We further evaluated the framework’s behavior across different LLM backbones in terms of F1-score, as shown in Figure 3. The results indicate that recent state-of-the-art models uniformly excel at understanding and identifying PFSMs for lightweight protocols. However, for large-scale industrial protocols,

models such as Grok and ChatGPT demonstrate superior capability in extracting structured state machines compared to others. Despite the strong performance of these backbones, the structural challenges described above (e.g., the branching in EDHOC) (Selander et al., 2021) still pose significant hurdles, necessitating the specialized mechanisms within the LLM-PF framework to achieve viable results. Additionally, incomplete operational descriptions, such as the omitted SQN handling mechanism in 5G AKA (p.35) (3GPP, 2022), remain a bottleneck.

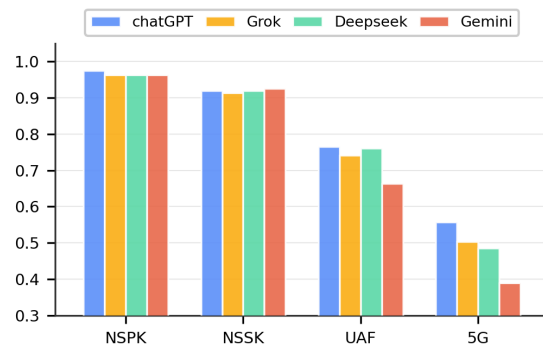


Figure 3: Comparison of F1-scores for global state machine extraction across different large language models (LLMs) and protocols.

5.3.2 Baseline Comparison (RQ2)

To evaluate the effectiveness of our framework, we established a strong baseline using direct prompt engineering: we instructed the LLM to generate the complete PFSM in a single pass using the raw protocol documents as input. Both the baseline and our LLM-PF framework were evaluated using the same automated metrics.

As shown in Figure 4, the comparative analysis

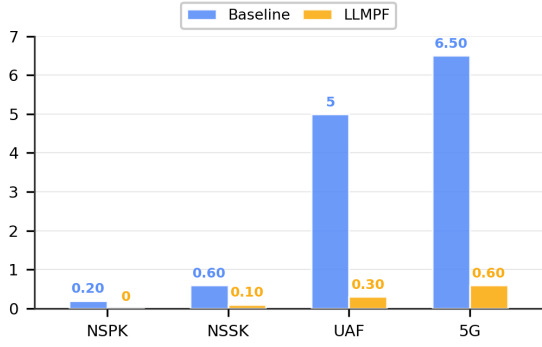


Figure 4: Comparison of average structural errors during global state machine extraction between LLM-PF and Baseline across different protocols.

reveals a clear dependency on protocol scale. For small-scale protocols (e.g., NSSK), the baseline demonstrates acceptable competence, and LLM-PF achieves only marginal improvements. This is primarily because the rectification module in LLM-PF corrects minor edge inconsistencies that are relatively sparse in simple protocols. However, in the context of large-scale protocols (e.g., 5G AKA), the baseline suffers a catastrophic drop in performance. This is due to the complexity of the PFSM structure and the difficulty of capturing inter-state dependencies and internal operations within the context window of a single forward pass. In this scenario, the advantage of LLM-PF’s hierarchical design becomes significant. Notably, for protocols with extensive documentation, the baseline often exceeds the context window limits of the LLM, resulting in incomplete or null generations. In contrast, LLM-PF’s modular design effectively manages context length and structural complexity, ensuring robust generation across all scales.

5.3.3 System Analysis (RQ3)

We conducted an ablation study to investigate the specific contributions of the global PFSM extraction and the rectification mechanism.

The Necessity of the Global PFSM. While the inability of LLMs to accurately identify all states in a single pass necessitates a decomposed strategy, a naive ‘Chain Extraction’ approach—extracting states sequentially based on predecessors—remains problematic. Although this method bypasses immediate context window limits, our analysis suggests that LLMs struggle to maintain coherence across multiple turns of dialogue. Specifically, such a sequential method fails to capture the migratory

relationships between states effectively and suffers from the progressive accumulation of hallucinations along the extraction chain. Consequently, we prioritized the Global PFSM extraction, which enables the model to rapidly capture the overarching state transition logic and message interdependencies at once, providing a robust structural guide for subsequent detailed extraction.

The Impact of the Rectification Mechanism. To quantify the impact of the rectification module, we conducted 10 independent trials of global PFSM extraction and measured the average number of structural errors, as shown in Figure 4. For small-scale protocols, direct extraction proves relatively reliable (averaging 0.2 errors), and our module successfully eliminates these residual inaccuracies. In contrast, large-scale protocols exhibit instability in direct extraction, yielding an average of 5-7 structural errors per run. The introduction of the rectification module significantly mitigates these issues, reducing the average unresolved errors to just 0.3 - 0.6. This demonstrates that the feedback loop is essential for ensuring the structural validity of complex protocol models.

6 Conclusion

This paper investigates the application of LLMs to cryptographic protocol flow extraction. We propose a standardized protocol description that can be comprehensively transformed into formal models and sequence diagrams. Leveraging LLMs, we introduce a hierarchical extraction and rectification algorithm framework that decomposes the task, thereby enhancing the capability to extract large-scale protocols. Experimental results on 10 protocols of varying complexity and four mainstream LLMs indicate that although LLMs are feasible, their performance on complex flows remains sub-optimal. Specifically, our framework significantly improves the capability of extracting cryptographic protocols compared to the baseline generated by directly using LLMs; however, we acknowledge that current limitations, such as the inability to recognize document cross-references, still exist. Future work will focus on addressing these issues, extending the evaluation to a broader range of protocols, and conducting security analysis based on the standardized descriptions.

649 Limitations

650 While our proposed framework demonstrates
651 promising capabilities, its effectiveness is currently
652 influenced by several factors. The reliance on large
653 models presents practical challenges, including sig-
654 nificant time and computational costs, and the ap-
655 proach has thus far been validated on a limited set
656 of protocol types and scales, with extensive test-
657 ing across diverse models and protocol categories
658 remaining for future work.

659 Furthermore, while our framework mitigates hal-
660 lucinations to a certain extent, it cannot entirely
661 eliminate the LLM from drawing erroneous con-
662 clusions. We also observe that ensuring high sta-
663 bility of model outputs can be difficult; even with
664 deterministic settings and our content completion
665 mechanism—which mitigates many common er-
666 rors—the system may not fully address rare, ex-
667 treme failures. Additionally, the quality of the gen-
668 erated PFSM is heavily dependent on the clarity of
669 the standard documentation. Omissions in proto-
670 col descriptions or dependencies that span multiple
671 documents can significantly impair the accuracy of
672 the generation results.

673 To address these inherent uncertainties and up-
674 hold the highest standards of reliability, we have in-
675 corporated a human-in-the-loop strategy. Although
676 LLM-PF strives to minimize human intervention,
677 the ultimate goal of formal verification requires
678 bridging the gap between symbolic models and real-
679 world scenarios. Consequently, certain details still
680 necessitate expert determination to ensure the for-
681 mal model accurately reflects protocol semantics,
682 such as the identification of principal confirmations
683 discussed in Section 5.3. This expert validation and
684 correction are crucial to guarantee the integrity of
685 the final formal specifications.

686 References

687 3GPP. 2022. [3GPP TS 33.501: Security architecture](#)
688 [and procedures for 5g system](#). Technical Specifica-
689 [tion 3GPP TS 33.501](#). Release 17.

690 Martín Abadi, Bruno Blanchet, and Cédric Fournet.
691 2017. The applied pi calculus: Mobile values, new
692 names, and secure communication. *Journal of the*
693 *ACM (JACM)*, 65(1):1–41.

694 Abdullah Al Ishtiaq, Sarkar Snigdha Sarathi Das, Syed
695 Md Mukit Rashid, Ali Ranjbar, Kai Tu, Tianwei Wu,
696 Zhezheng Song, Weixuan Wang, Mujtahid Akon, Rui
697 Zhang, and 1 others. 2024. Hermes: Unlocking secu-
698 rity analysis of cellular network protocols by syn-

699 thesizing finite state machines from natural language
700 specifications. In *33rd USENIX Security Symposium*
701 *(USENIX Security 24)*, pages 4445–4462.

Luca Arnaboldi and Roberto Metere. 2019. Towards
702 a data centric approach for the design and verifica-
703 tion of cryptographic protocols. *arXiv preprint*
704 *arXiv:1910.02656*. 705

David Basin, Michel Keller, Saša Radomirović, and
706 Ralf Sasse. 2015. Alice and bob meet equational the-
707 ories. In *Logic, Rewriting, and Concurrency: Essays*
708 *Dedicated to José Meseguer on the Occasion of His*
709 *65th Birthday*, pages 160–180. Springer. 710

Bruno Blanchet. 2012. Security protocol verification:
711 Symbolic and computational models. In *Internation-*
712 *al conference on principles of security and trust*,
713 pages 3–29. Springer. 714

Michael Burrows, Martin Abadi, and Roger Needham.
715 1990. A logic of authentication. *ACM Transactions*
716 *on Computer Systems (TOCS)*, 8(1):18–36. 717

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann,
718 Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-
719 cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and
720 1 others. 2025. Gemini 2.5: Pushing the frontier with
721 advanced reasoning, multimodality, long context, and
722 next generation agentic capabilities. *arXiv preprint*
723 *arXiv:2507.06261*. 724

Cristian Curaba, D’Ambrosi Denis, and Alessandro
725 Minisini. 2024. Cryptoformaleval: Integrating large
726 language models and formal verification for auto-
727 mated cryptographic protocol vulnerability detection.
728 In *The First Workshop on System-2 Reasoning at*
729 *Scale, NeurIPS’24*. 730

Martin Duclos, Ivan A Fernandez, Kaneesha Moore,
731 Sudip Mittal, and Edward Ziegler. 2024. Utiliz-
732 ing large language models to translate rfc protocol
733 specifications to cpsa definitions. *arXiv preprint*
734 *arXiv:2402.00890*. 735

FIDO Alliance. 2017. [Fido uaf protocol specification](#).
736 Proposed standard, FIDO Alliance. Accessed: 2026-
737 01-04. 738

Oded Goldreich. 2003. Cryptography and cryp-
739 tographic protocols. *Distributed Computing*,
740 16(2):177–199. 741

Young-Hoon Goo, Kyu-Seok Shim, Min-Seob Lee, and
742 Myung-Sup Kim. 2019. Protocol specification ex-
743 traction based on contiguous sequential pattern algo-
744 rithm. *IEEE Access*, 7:36057–36074. 745

Mohammad Mazyad Hazzazi, Raja Rao Budaraju, Zaid
746 Bassfar, Ashwag Albakri, and Sanjay Mishra. 2023.
747 A finite state machine-based improved cryptographic
748 technique. *Mathematics*, 11(10):2225. 749

Amey Hengle, Prasoon Bajpai, Soham Dan, and Tan-
750 moy Chakraborty. 2025. Multilingual needle in a
751 haystack.

858 A.2 Abstraction

859 Each entity PFSM is abstracted as

$$860 \text{FSM}_e \triangleq (S_e, s_e^0, O_e, \theta_e), \quad (7)$$

861 where S_e is the set of states, s_e^0 is the initial state,
862 $O_e(s)$ returns the onEntry operation sequence at
863 state s , and $\theta_e(s)$ returns the outgoing transitions
864 at s (e.g., AUTO/CONDITIONAL).

865 **Terms, patterns, and substitution.** Terms and
866 patterns are tuples built from variables and con-
867 stants:

$$868 t ::= x \mid c \mid (t_1, \dots, t_n),$$

$$869 p ::= x \mid c \mid (p_1, \dots, p_n).$$

870 A substitution σ maps variables to terms; applying
871 σ to t is written $t\sigma$. We write $\text{match}(p, t\sigma) = \sigma'$
872 when p matches $t\sigma$ and yields bindings σ' , other-
873 wise $\text{match}(p, t\sigma) = \perp$.

874 A.3 Event Language and Event Programs

875 **Events.** We abstract onEntry and transition ac-
876 tions into events:

$$877 \alpha \in \text{Ev} ::= \text{SEND}(ch, t) \mid \text{RCV}(ch, p)$$

$$878 \mid \text{SPLIT}(t \Rightarrow \vec{x}) \mid \text{OP}(x := f(\vec{t})).$$

879 Fresh-name generation is modeled as a dedicated
880 event $\text{RAND_GEN}(x)$, which compiles to the
881 ProVerif statement new . We use vector notation \vec{x}
882 and \vec{t} to denote tuples (e.g., $\vec{x} = (x_1, \dots, x_k)$ and
883 $\vec{t} = (t_1, \dots, t_n)$), and write $\vec{t}\sigma$ for $(t_1\sigma, \dots, t_n\sigma)$.

884 **Event programs.** Role behavior is expressed as
885 an event program:

$$886 P ::= 0 \mid \alpha; P \mid \text{if } \varphi \text{ then } P_1 \text{ else } P_2$$

$$887 \mid \text{case } t \text{ of } \{p_i \Rightarrow P_i\}_{i=1}^k$$

$$888 \text{ default } \Rightarrow P_0,$$

889 where φ is a boolean guard (from PFSM condi-
890 tionals) and **case** uniformly captures “receive-then-
891 discriminate” branching (e.g., by tag/shape).

892 A.4 Operational Semantics (Short Form)

893 **Configuration and multiset scheduling.** A con-
894 figuration is

$$895 C \triangleq \langle \mathcal{N}, \mathcal{I}, \sigma \rangle,$$

896 where \mathcal{N} is the set of generated fresh names, \mathcal{I} is a
897 *multiset* of concurrently running instances, and σ
898 is the current environment/substitution.

899 An instance is written $inst \triangleq \{\alpha; P\}$, mean-
900 ing “execute event α then continue with P ”. We
901 use multiset union \uplus , multiset difference \setminus_m , and
902 multiset inclusion \subseteq_m .

903 We use the following multiset-update notation
904 (replace one instance by its continuation):

$$905 \mathcal{I}[inst \mapsto P] \triangleq (\mathcal{I} \setminus_m \{inst\}) \uplus \{P\}.$$

906 **Step relation.** Each semantic step emits a *label*
907 ℓ that is later mapped to ProVerif code:

$$908 \langle \mathcal{N}, \mathcal{I}, \sigma \rangle \xrightarrow{\ell} \langle \mathcal{N}', \mathcal{I}', \sigma' \rangle.$$

909 **Core rules.**

$$\frac{inst = \{\text{RAND_GEN}(x); P\} \uplus \mathcal{I} \quad n \notin \mathcal{N}}{\langle \mathcal{N}, \mathcal{I}, \sigma \rangle \xrightarrow{\text{new}(x)} \langle \mathcal{N} \cup \{n\}, \mathcal{I}[inst \mapsto P], \sigma[x \mapsto n] \rangle}$$

$$\frac{inst = \{\text{SEND}(ch, m); P\} \uplus \mathcal{I}}{\langle \mathcal{N}, \mathcal{I}, \sigma \rangle \xrightarrow{\text{out}(ch, m\sigma)} \langle \mathcal{N}, \mathcal{I}[inst \mapsto P], \sigma \rangle}$$

$$\frac{inst = \{\text{RCV}(ch, p); P\} \uplus \mathcal{I} \quad \text{match}(p, v) = \sigma'}{\langle \mathcal{N}, \mathcal{I}, \sigma \rangle \xrightarrow{\text{in}(ch, p)} \langle \mathcal{N}, \mathcal{I}[inst \mapsto P], \sigma \cup \sigma' \rangle}$$

$$\frac{inst = \{\text{SPLIT}(t \Rightarrow \vec{x}); P\} \uplus \mathcal{I} \quad \text{match}(\vec{x}, t\sigma) = \sigma'}{\langle \mathcal{N}, \mathcal{I}, \sigma \rangle \xrightarrow{\text{let } \vec{x} = t\sigma} \langle \mathcal{N}, \mathcal{I}[inst \mapsto P], \sigma \cup \sigma' \rangle}$$

$$\frac{inst = \{\text{OP}(x := f(\vec{t})); P\} \uplus \mathcal{I}}{\langle \mathcal{N}, \mathcal{I}, \sigma \rangle \xrightarrow{\text{let } x = f(\vec{t}\sigma)} \langle \mathcal{N}, \mathcal{I}[inst \mapsto P], \sigma[x \mapsto f(\vec{t}\sigma)] \rangle}$$

909 A.5 Label-to-ProVerif Mapping (Code 910 Generation)

911 We map each transition label ℓ to ProVerif code via
912 $\Phi(\ell)$. We write $(\vec{x} : \vec{\tau})$ to denote $(x_1 : \tau_1, \dots, x_k : \tau_k)$,
913 where the types are declared in \mathcal{H} . The map-
914 ping Φ from labels to ProVerif statements is:

$$915 \Phi(\text{new}(x)) \triangleq \text{new } x : \text{bitstring};$$

$$\Phi(\text{out}(ch, m)) \triangleq \text{out}(ch, m);$$

$$\Phi(\text{in}(ch, p)) \triangleq \text{in}(ch, p);$$

$$\Phi(\text{let } \vec{x} = t) \triangleq \text{let } (\vec{x} : \vec{\tau}) = t \quad (8)$$

$$\text{in}$$

$$\Phi(\text{let } x = f(\vec{t})) \triangleq \text{let } x : \tau = f(\vec{t})$$

$$\text{in.}$$

916 **A.6 From FSM to Event Programs**
 917 **(Control-Flow Extraction)**

918 We compile a state s of entity e into an event pro-
 919 gram by concatenating its onEntry code and its
 920 outgoing transitions:

$$\text{CompileState}(e, s) \triangleq \text{CompileOps}(O_e(s));$$

$$\text{CompileTrans}(e, s, \theta_e(s)). \quad (9)$$

921 The role-level event program starts from the initial
 922 state:

$$923 p_e \triangleq \text{CompileState}(e, s_e^0).$$

924 **A.7 End-to-End Summary**

925 The complete translation can be summarized as:

926 **Step 1. Extract PFSM.**

$$927 \mathcal{P} \mapsto (\mathcal{E}, (\text{FSM}_e)_{e \in \mathcal{E}}).$$

928 **Step 2. Compile to event programs.**

$$929 p_e \triangleq \text{CompileState}(e, s_e^0).$$

930 **Step 3. Run semantics (labels).**

$$931 p_e \Downarrow \ell_1 \cdots \ell_n.$$

932 **Step 4. Emit ProVerif code.**

$$933 \text{Proc}_e \triangleq \Phi(\ell_1); \cdots ; \Phi(\ell_n).$$

934 **Step 5. Assemble output.**

$$935 \text{PV}(\mathcal{P}) = \mathcal{H} \cup \bigcup_{e \in \mathcal{E}} \text{Proc}_e \cup \text{Proc}_{\text{init}}.$$

936 **A.8 Notation (Quick Reference)**

937 \mathcal{P} Input protocol PFSM.

938 \mathcal{E} Entity set; e denotes an entity/role.

939 FSM_e
 940 Local PFSM $(S_e, s_e^0, O_e, \theta_e)$ of entity e .

941 Proc_e
 942 ProVerif process for e ; $\text{Proc}_{\text{init}}$ is the initializa-
 943 tion process.

944 t, p
 945 Term and pattern; σ is a substitution/environment;
 946 $t\sigma$ denotes substitution application.

947 C Configuration $C = \langle \mathcal{N}, \mathcal{I}, \sigma \rangle$, where \mathcal{N} is the
 948 fresh-name set and \mathcal{I} is a multiset of instances.
 949

$inst$ 950

Instance $inst = \{\alpha; P\}$ (execute event α then
 951 continue with P). 952

ℓ, Φ 953

Transition label ℓ and its label-to-ProVerif map-
 954 ping Φ . 955

956 **A.9 An Example**

957 Figure 5 illustrates a concrete instantiation of
 958 PFSM on a simple cryptographic protocol in which
 959 Alice sends a digitally signed value to Bob, who
 960 verifies the signature. The left-hand side of the fig-
 961 ure shows the PFSM representation of the protocol,
 962 while the right-hand side presents its corresponding
 963 formal realization.

964 At the protocol level, the participating entities
 965 Alice and Bob correspond to elements in the en-
 966 tity set \mathcal{E} , and their initial knowledge is specified
 967 by the knowledge distribution function Ks . These
 968 components define the static context under which
 969 protocol execution begins.

970 The behavior of each participant is modeled as a
 971 finite state machine. As shown in Figure 5, Alice
 972 is represented by a single state unit whose action
 973 set includes a composite action that computes a
 974 signature and sends the resulting message to Bob.
 975 Bob is similarly modeled by a single state unit,
 976 whose actions include receiving Alice's message
 977 and performing signature verification. These ac-
 978 tions correspond to the action sets $\theta_{e,s}$ defined in
 979 PFSM.

980 Protocol control flow is captured by the transi-
 981 tion functions $\theta_{e,s}$. In this example, Bob's transi-
 982 tion is conditional: successful verification leads to
 983 a success state, while failure leads to an alternative
 984 outcome. This illustrates how PFSM explicitly rep-
 985 resents conditional branching and protocol flow at
 986 the state level.

987 Overall, Figure 5 demonstrates how PFSM or-
 988 ganizes protocol entities, knowledge, actions, and
 989 control flow into a structured intermediate repre-
 990 sentation. Beyond improving interpretability, this
 991 design also facilitates the systematic generation of
 992 formal code for verification tools such as ProVerif,
 993 by aligning PFSM constructs with the operational
 994 abstractions employed in formal verification frame-
 995 works.

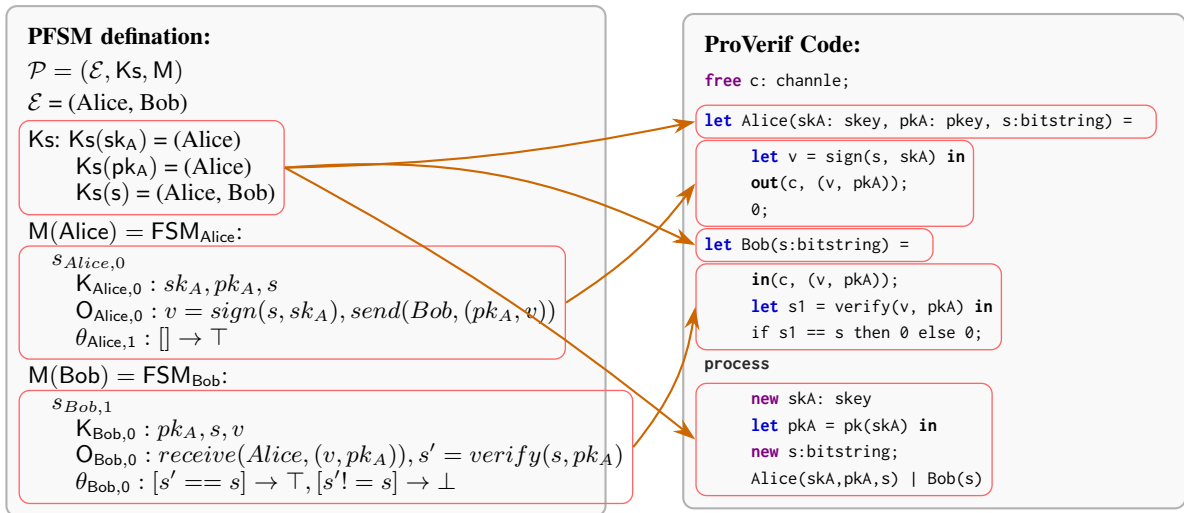


Figure 5: Mapping PFSM definition to Formal Code