



NSR Modernization

David Brown, Ben Shu, Ramon Arcilla, Boris Pritychenko National Nuclear Data Center, Brookhaven National Laboratory, Upton NY, USA

NRDC, Madrid Spain, 19 June 2025



NSR Workflow





Previous Workflow



- All steps handled by human beings
 - But we have no more contractors as of April 2025



Previous Workflow



- Blue steps are rate determining; are the target for improvements
 - Legacy process must be carefully documented



Web scraping

• Article citations formerly obtained by people searching in websites

• Ex. Physical Review C (PRC)

- Click on one article, download BibTeX
- Rinse and repeat for an entire issue
- Vol 11, Issue 3 85 articles added
- This can be automated by web scraping tools

ł	Export citation
E	Export citation
(Choose format for download:
0	BibTeX C EndNote (RIS) Download Citation
	<pre>@article{PhysRevC.111.044301, title = {Angular momentum geometry of pseudospin partner bands in \$^{133}\mathrm{Cs} author = {Chen, Yeruoxi and Chen, Q. B.}, journal = {Phys. Rev. C}, volume = {111}, issue = {4}, pages = {044301}, numpages = {11},</pre>
	year = {2025}, month = {Apr}, publisher = {American Physical Society},



Web scraping

Article citations formerly obtained by people sea We now have RIS & Bibtex translators and web scrapers Ex. Physical for several journals d Citation Click on one ry of pseudospin partner bands in \$^{133}\mathrm{Cs}\$} Rinse and r Q. B.}, • Vol 11, Issu (thank you Ben and Boris!) year = {2025} This can be automated by month = {Apr}. publisher = {American Physical Society} (10.1102/DhueDouC.111.04420) web scraping tools



Slightly Improved Workflow



- Simple scrapers have sped up ingest from 1 week to 1 hour
- We've sped up article ingest, but we can go faster!



Development Workflow

scrape-from-prc

- Automatically scrapes and processes PRC articles
- Connects to PRC website to obtain articles
- Downloads BibTeX, converts to NSR "exchange"

<pre>@article{PhysRevC.111.044301,</pre>	
title = {Angular momentum geometry of pseudospin partner bands in 133 author = {Chen Veruoxi and Chen O B }	atnrm{CS}\$};
journal = {Phys. Rev. C},	
volume = {111},	
issue = {4},	
pages = {044301},	<keyno>2025CHP</keyno>
numpages = {11},	<hr/> thistory >C20250409
year = {2025},	
month = {Apr};	<pre><coden>JOUR PRVCA 111 044301</coden></pre>
doi = {10.1103/PhysRevC.111.044301},	<refrence>Phys.Rev. C 111, 044301 (2025)</refrence>
<pre>url = {<u>https://link.aps.org/doi/10.1103/PhysRevC.111.044301</u>}</pre>	<authors>Y.Chen, Q.B.Chen</authors>
}	<title>Angular momentum geometry of pseudospin partner bands in {+133}Cs</title>
	<doi>10.1103/PhysRevC.111.044301</doi>



update-prc-errata

- Checks for updates to existing NSR entries
- Preserves key number, title, keywords, selectors, etc.





save-sql-backup

- Saves a copy of NSR using mariadb-dump
- Runs before any database changes
- Can be used as a restore point, just in case



nsrprep

- Based on legacy web app used by compilers
 - Performs format/validation checks
 - Looks for unknown author names
- Generates selectors used for NSR website searching
 - Only generates these for **keyworded** entries
 - If an entry is not keyworded, it cannot be found with:
 - Quick Search: Nuclide, Reaction
 - Indexed Search: Nuclide, Subject, Reaction, Z, A



update-database

- Makes changes to nsr_automated
 - Internal only not deployed to NSR website
- Inserts new entries for standard PRC articles
- Updates existing entries for PRC errata
- **Issue:** Currently no way to fix errors before insertion
 - i.e. format errors, PRC typos
 - Fixes require direct SQL changes or restarting pipeline



Passed Cat Dunn crea	ted pipeline for commit 23ff2459	17 hours ago, finished 17 hours ago			
r main					
est branch CO 6 jobs (7 minutes 3 seconds, queued for 2 s	seconds			
vipeline Jobs 6 T	Fests 0				
build	collect	update-errata	save-sql-backup	generate-selec	update-database

- Pipeline functional as of April 3rd, 2025
 - No setup/installation only GitLab access required
 - Click a few buttons, go do something else



A Much Faster Workflow



- Full pipeline takes minutes to run
- We are not taking advantage of the full power of gitlab



NSR Exchange format can be modernized

The current NSR "exchange" format...

<keyno > 2023wizz <HISTORY > A20230605 <CODEN > CONF Sacramento... <REFRENCE> Proc.15th.Intern... Citation information <AUTHORS > D.Wiarda, ... <TITLE > Modernization... <KEYWORDS> [N/A] <SELECTRS> [N/A]

Unique ID Entry/modification dates Source type and information List of authors Formatted title Sentences about content Structured search keys



NSR Exchange format can be modernized

• ...can be translated into JSON...





NSR Exchange format can be modernized

- ...but the exchange files don't tell the full story
- NSR's SQL database has supporting tables:





What does this get us?

A JSON format enables:

- A JSON-schema
- Format and correctness checkers
- An editor
- Better NSR archive option
- An AI/ML-friendly format!

Even without these, with gitlab we have

- Version control
- Built in review step (branch merge process)
- We use this to great effect in ENDF & JEFF





Nuclear Science References (NSR)

Nuclear Science literature database, containing about 250,000 articles covering more than 120 years of nuclear science research. Approximately 80 journals are routinely scanned for relevant articles; it also includes PhD thesis and private communications. Each article in NSR is assigned a distinctive 8-character primary key, known as 'key number', and a set of keywords that briefly describe the article's content. These keywords follow strict rules since they are used in web search forms. Typically, an article is incorporated in NSR within a few weeks after its web publication. About 4,000 articles are added per year. All the compilation effort is supported by the US Nuclear Data Program and coordinated by the National Nuclear Data Center at Brochkaven National Laboratory, which is also responsible for its web dissemination.



242,120 Articles	196,773 Keyworded Abstracts	108,272 Authors	8,578 Ruclides	8,914 Reactions	382 Decays
Deposition Summary				Dataset Details	
Depositor:	В	oris Pritychenko		Total articles:	242120
Contact:	р	ritychenko@bnl.gov		Total keyworded	196773
Deposition date:	2	022/04/06		abstracts.	400000
Last modified:	2	022/04/06		Total authors:	108272
DOI:	1	0.18139/nndc.nsr/1866	Nuclides:	8578	
			Reactions:	8914	
Latest Dataset 🗟				Decays:	382
File	Date	•			

https://www.nndc.bnl.gov/nsrarchivals/

i 0	+ 😔	📼 ENDF / 💌 libra	ry / 😨 n	utrons / M	erge requ	ests / 11	110					
D 4 11	 2	Draft: Re	view	/n-01	7_CI	_03	7				Edit	ode ~
et occircit of go t		33 Open Gusta	vo Nobre	requested	to merge	Review	/n-017_0	1_037 🛱 into	phase2 🕃 2 m	onths ago		
Project		Overview 0	Commi	r 1 Di	nalinas	1 CF	anger				Add	a to-do item
😵 neutrons		Over view 0	Commi		peunes		anges	5			Au	a to do item
🖈 Pinned	~	Phase I Re	view								0 Assignees	Edit
Issues	208										None - assign yourset	
	20	ENDF/B-VII	I.1, Neu	trons S	ublibr	ary					0 Peviewers	Edit
Merge requests	20	- Eilenamou	017 CL (27 andf							Approval is optional Assign	Lun
සී Manage	>	 Sublibrary 	Neutron:	S7.enui							None - assign vourself	
岗 Plan	>	Material: 1	7-Cl- 37 (MAT=1731)								
2 Flain		 Evaluators 	Sayer,Gu	ber,Kawar	io,Hanse	lman+					Labels	Edit
/> Code	~	 Submitter: 	Kenneth I	lanselman							None	
Merge requests	25	 Submitter 	email: kha	nselman@	lanl.gov							
mergerequeets		 Review for 	m genera	ted at: 11/0	4/2025	13:43:38					Milestone	Edit
Repository		 Reviewer: 									None	
Branches		Reviewer e	mail:									
Commite		• Date:				-					Time tracking	ф +
Tags		Instructions									No estimate or time spent	
Repository graph		You are being as	ked to rev	iew the ev	aluation	for n-01	7_C1_03	7.endf from b	ranch phase1.P	Please edit this file on the	5 Dentisia ant	
		Review/n-017_0	1_037 me	rge reque	st tracke	r page at	11110. Y	ou can do this :	imply by clickin	g the "edit" button and filling the empty fields		
Compare revisions		This document u	ses the N	arkdown f	ormat (s	ee https:	//docs.g	itlab.com/ee/u	ser/markdown.h	tml for a reference.).	3	
Snippets		File Contents										
Locked files		File Contents										
0		Questions to consider with each dataset: Is the current dataset an improvement over the existing dataset in the ENDF evaluation? Is the										
8. Brilla	,	current dataset	current dataset (e.g. cross section) some form of standard or reference? If so, how does the dataset compare to the reference? Does the									
D Secure	>	current set take into account all relevant differential data? Data and covariance adequacy should be given on a scale of 0-5 with 5 being the bioheet Llex NA if not driven in the KNDF file.										
Deploy	>	ingricus. Obe										
		Reaction	1 2	3 4	6	32	33	Data	Cov.	Comments		
6*11-1-		(MT)						Adequacy	Adequacy			



With Gitlab, we can ensure quality while getting speed



 With a QA process in hand, and an automated workflow, now we can target the hardest step



Keywording

Currently a compiler must produce keyworded abstract – time consuming and requires somewhat special skills

As of FY25, NNDC has no contractors

Options:

- Teach more people to do it (ENSDF/ENDF evaluators, students, NNDC staff?)
- Automate it?



Nuclear Science References Coding Manual

D.F. Winchell National Nuclear Data Center Brookhaven National Laboratory Upton, New York, USA

May 24, 2007

https://www.nndc.bnl.gov/nsr/docs/nsr-coding-manual.pdf



Future Workflow



- Need to fill in missing steps
 - Keyword assignment via chatNSR
 - Tools for human reviewers to make edits



chatNSR: An AI-Enhanced Nuclear Science References Knowledge Base

Ramon Arcilla



Motivation

Why automate keyword abstract generation? **Volume Challenge**

• Growing nuclear science literature (~5,000+ annually) strains current manual keywording capacity.

Time Intensive

• Manual keywording requires ~3.5 hours (average) per article.

Limited Coverage

Manual keywording limits timeliness and breadth of NSR knowledge base content.

Discovery Issues

- Inconsistent keywording limits paper discoverability.
 Cost Inefficiency
- Manual keywording demands significant resources.







The chatNSR Project

Current Focus:

• Keywording, the most tedious and costly activity

Key AI Technologies in Use

- Open-Source Large Language Models (LLMs)
 - Understand user queries/instructions and generate keyword abstract
- Cache-Augmented Generation (CAG)
 - Loads entire article into LLM's extended context window (Knowledge Cache) for efficiency, more accurate keywording
- Retrieval-Augmented Generation (RAG)
 - Finds contextually relevant parts of article in vector DB for LLM's use in keywording
- System Prompt Engineering (SYSPROMPT)
 - User instructions to guide LLMs for consistent, accurate keywording



The chatNSR Project

Architectural Overview: Two-Path Automation Workflow

- Document Input
 - Article enters the automation pipeline
- Size Assessment
 - System determines article's token count
- Path Selection
 - CAG for smaller articles (≤100k tokens = 100 pages), RAG for larger ones (>100k tokens)
- Keyword Output
 - Keyword abstracts generated via optimal path



Cache-Augmented Generation (CAG)

Workflow

- PDF Ingestion
 - Article enters system and undergoes initial parsing
- Markdown Conversion
 - Content structure preserved in clean Markdown format
- Context Preloading
 - Full article loaded with Knowledge Cache creation
- System Prompt Processing
 - LLM guided by engineered instructions
- Response Output
 - Keyword abstracts delivered to user



Retrieval-Augmented Generation (RAG)

Workflow

- Document Processing
 - PDF ingestion and conversion to markdown format.
- Chunking & Embedding
 - Document segmentation and vector representation creation.
- Retrieval Mechanism
 - Query-based selection of relevant document sections.
- Generation with Prompting
 - LLM processes retrieved chunks with specialized instructions.
- Output Delivery
 - Keyword abstract returned to requesting user.



Advantages of the Two-Path Automation Approach

Universal Scalability

- System handles articles of any length.
- Processing adapts to content size automatically.

Consistent Quality

- Keyword abstract coherence maintained across both workflows.
- Keyword abstract remains relevant regardless of article size.

Resource Efficiency

- Computational resources allocated based on article needs.
- Simpler approach (CAG) used when possible.

Workflow Integration

- Unified API endpoint despite dual backend systems.
- Seamless integration with existing document processes.



EXAMPLE: Keywording D.A. Brown's article "Zirconium Evaluations for ENDF/ B-VII.2 for the Fast Region"



Human-Generated

COMPILATION ^{90, 91, 92, 93, 94, 95, 96}Zr(n, x), E=0=20 MeV; calculated, evaluated σ , $\sigma(\theta)$ including resonance region using EMPIRE-3.1 code. Compared with available data and evaluated databases.

chatNSR-Generated

COMPILATION ⁹⁰ ⁹¹ ⁹² ⁹³ ⁹⁴ ⁹⁵ ⁹⁶Zr(n,x), E=0-20 MeV; calculated, evaluated σ , $\sigma(\theta)$ including resonance region using EMPIRE-3.1 code. Compared with available experimental data and evaluated databases (ENDF/B-VII.1, JENDL-4. 0).

Current Status

Testing

- Dozen open-source LLM's tested
- Dozens of articles (different topics) for LLM training
- Perplexity Al's *r1-1776:70b* LLM has highest accuracy

Accuracy

- 80-85% accuracy for articles \leq 10 pages
- More work needed for articles > 10 pages

Methods and Tools in Use

- Continuous human feedback to LLM
- Al-assisted coding and innovative prompt engineering
- PhysBERT (LBNL) model for embedding and retrieval
- IBM Docling software for PDF-to-Markdown conversion
- Perplexity Al's *r1-1776* LLM for keyword generation



Summary and Conclusion

Keywording Automation

- Continuous human feedback improves LLM's accuracy.
- Article (varied topics) high-volume ingestion and training enhances LLM's performance.
- CAG, RAG, System Prompts minimize/prevent LLM hallucination.
- 90% time reduction in keyword abstract generation
- PDF-to-Markdown conversion improves LLM's article understanding

On-premise Keywording

- Open-source LLMs and AI tools demonstrate cost-free, flexible, adaptable keywording
- More secure environment behind BNL FireWall
- Compliant with DOE rules and copyright laws



Future Direction

LLM capability and size

 Move to more advanced LLMs with > 70 billion parameters for higher accuracy.

Article Ingestion

• Massive article ingestion with one line command: >>> genabs articles-directory/ keywords-directory/

Multi-modal support

Processing of article's tables, figures, and equations

Computational resources

• Upgrade of computational resources to meet increased processing needs.



If we can (nearly) fully automate NSR, then both XUNDL and EXFOR should be possible





Thank you so much for your attention!

This work is sponsored by the Office of Nuclear Physics, Office of Science of the U.S. Department of Energy under Contract No. DE-SC0012704 with Brookhaven Science Associates, LLC.

35