Rigorous Error Analysis for Logarithmic Number Systems

Presented at the 32nd IEEE International Symposium on Computer Arithmetic ARITH 2025

Thanh Son Nguyen

University of Utah

thahnson@cs.utah.edu

Alexey Solovyev

University of Utah solovyev.alexey@gmail.com

Ganesh Gopalakrishnan

University of Utah ganesh@cs.utah.edu

Mark G. Arnold Lehigh University maab@lehigh.edu





www.xlnsresearch.com github.com/xlnsresearch



https://github.com/soarlab/RigorousErrorLNS



www.xlnsresearch.com website (also xlnsresearch.github.io) = Logarithmic Number Systems (LNS) papers



www.xlnsresearch.com website (also xlnsresearch.github.io) = Logarithmic Number Systems (LNS) papers

github.com/xInsresearch open-source repositories for LNS

xins configurable Python library for LNS tensors (analogous to Numpy)
 available as PyPi wheel via pip install xlns
 play with it in Python import xlns as xl



www.xlnsresearch.com website (also xlnsresearch.github.io) = Logarithmic Number Systems (LNS) papers

github.com/xInsresearch open-source repositories for LNS

xins configurable Python library for LNS tensors (analogous to Numpy)
 available as PyPi wheel via pip install xlns
 play with it in Python import xlns as xl
 supports plug-in configurations for novel LNS arithmetic algorithms
 actively seeks open-source contribution of these
 algorithms used in this paper are available import xlnsconf.utah_tayco_ufunc



www.xlnsresearch.com website (also xlnsresearch.github.io) = Logarithmic Number Systems (LNS) papers

github.com/xInsresearch open-source repositories for LNS

xins configurable Python library for LNS tensors (analogous to Numpy)
 available as PyPi wheel via pip install xlns
 play with it in Python import xlns as xl
 supports plug-in configurations for novel LNS arithmetic algorithms
 actively seeks open-source contribution of these
 algorithms used in this paper are available import xlnsconf.utah_tayco_ufunc

```
>>> import xlns as xl
>>> xl.xlnsnp([1,2,3])+4
xlnsnp([xlns(5.00000016087236) xlns(5.999999933686084) xlns(7.000000011403832)])
```



www.xlnsresearch.com website (also xlnsresearch.github.io) = Logarithmic Number Systems (LNS) papers

github.com/xInsresearch open-source repositories for LNS

xins configurable Python library for LNS tensors (analogous to Numpy)
 available as PyPi wheel via pip install xlns
 play with it in Python import xlns as xl
 supports plug-in configurations for novel LNS arithmetic algorithms
 actively seeks open-source contribution of these
 algorithms used in this paper are available import xlnsconf.utah_tayco_ufunc

```
>>> import xlns as xl
>>> xl.xlnsnp([1,2,3])+4
xlnsnp([xlns(5.000(0016)87236) xlns(5.999999933686084) xlns(7.000000011403832)])
>>> import xlnsconf.utol_tayco_ufunc
>>> xl.xlnsnp([1,2,3])+1
xlnsnp([xlns(4.9999974)724449) xlns(5.999999933686084) xlns(6.999999432996774)]
>>> xl.xlnssetf(9)
>>> xl.xlnsnp([1 2,3])+4
xlnsnp([xlns(4.994403908756931) xlns(5.995933468164773) xlns(7.006013412127704)])
```

Advantages of LNS

Cheaper multiply, divide, square root Good for applications with high proportion of multiplications <u>log(3)</u>



$$\log_2(2^p \times 2^q) = \log_2 2^{p+q} = p+q$$
$$\log_2(2^p/2^q) = \log_2 2^{p-q} = p-q$$





Usually, base b = 2



Usually, base b = 2





Discrete change in distance causes wobble in relative precision

Continuous change in distance means constant relative precision

Floating Point





Commercial Interest in LNS

European Union:	European Logarithmic Microprocessor (ELM)
Motorola:	LNS chip for satellite network
Yamaha:	Music Synthesizer
Boeing:	Aircraft controls
Interactive Machines,Inc.:	IMI-500: Animation for Jay Jay the Jet Plane
Advanced Rendering Technologies:	Hardware Ray-Tracing Engine
Cambridge/Microsoft:	HTK Hidden Markov Model Toolkit
Univ. of Tokyo:	N-body Gravity Pipeline (GRAPE): Gordon Bell Prize
NVIDIA:	LNS-MADAM and novel LNS summation

LNS Wordsize for Applications



Biological Intelligence uses Logarithms

Weber–Fechner law: ... response to sensory stimulus (light, sound,...) proportional to the logarithm of the stimulus.

Mental numbers uses logarithmic scale.

"biological synapses...26 levels in a logarithmic number system, thus storing ... 5 bits"



J. E. Volk, B Parhami, "Number Representation and Arithmetic in the Human Brain", <u>https://web.ece.ucsb.edu/~parhami/pubs_folder/parh20-iemcon-arithmtic-human-brain-final.pdf</u>

T. M. Bartol, Jr., et al., "Nanoconnectomic upper bound on the variability of synaptic plasticity", *eLife*, 2015. 10.7554/eLife.10778.002

Number Systems for Deep Neural Network Architectures: A Survey

Ghada Alsuhli¹, Vasileios Sakellariou¹, Hani Saleh, *Senior Member, IEEE*,¹, Mahmoud Al-Qutayri, *Senior Member, IEEE*,¹, Baker Mohammad, *Senior Member, IEEE*,¹, and Thanos Stouraitis¹

- [47] M. G. Arnold, T. A. Bailey, J. J. Cupal, and M. D. Winkel, "On the cost effectiveness of logarithmic arithmetic for backpropagation training on simd processors," in *Proceedings of International Conference on Neural Networks (ICNN'97)*, vol. 2. IEEE, 1997, pp. 933–936.
- [48] B. Parhami, "Computing with logarithmic number system arithmetic: Implementation methods and performance benefits," *Computers & Electrical Engineering*, vol. 87, p. 106800, 2020.
- [49] D. Miyashita, E. H. Lee, and B. Murmann, "Convolutional neural networks using logarithmic data representation," *arXiv preprint arXiv*:1603.01025, 2016.
- [50] A. Sanyal, P. A. Beerel, and K. M. Chugg, "Neural network training with approximate logarithmic computations," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 3122–3126.
- [51] S. A. Alam, J. Garland, and D. Gregg, "Low-precision logarithmic number systems: Beyond base-2," ACM Transactions on Architecture and Code Optimization (TACO), vol. 18, no. 4, pp. 1–25, 2021.
- [52] I. Kouretas and V. Paliouras, "Logarithmic number system for deep learning," in *International Conference on Modern Circuits and Systems Technologies (MOCAST)*. IEEE, 2018, pp. 1–4.
- [53] H. Saadat, H. Bokhari, and S. Parameswaran, "Minimally biased multipliers for approximate integer and floating-point multiplication," *IEEE Transactions on Computer-Aided Design of Integrated Circuits* and Systems, vol. 37, no. 11, pp. 2623–2635, 2018.
- [54] J. N. Mitchell, "Computer multiplication and division using binary logarithms," *IRE Transactions on Electronic Computers*, no. 4, pp. 512– 517, 1962.

- [55] Z. Babić, A. Avramović, and P. Bulić, "An iterative logarithmic multiplier," *Microprocessors and Microsystems*, vol. 35, no. 1, pp. 23– 33, 2011.
- [56] M. S. Ansari, B. F. Cockburn, and J. Han, "A hardware-efficient logarithmic multiplier with improved accuracy," in *Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2019, pp. 928–931.
- [57] R. Pilipović and P. Bulić, "On the design of logarithmic multiplier using radix-4 Booth encoding," *IEEE access*, vol. 8, pp. 64 578–64 590, 2020.
 [58] L. Harsha, B. R. Jammu, N. Bodasinei, S. Veeramachaneni, and N. M.
- [56] E. Hansha, B. K. Fallmur, N. Bodasingi, S. veccantachaneur, and N. M. SK, "A low error, hardware efficient logarithmic multiplier," *Circuits, Systems, and Signal Processing*, vol. 41, no. 1, pp. 485–513, 2022.
- [59] M. S. Kim, A. A. Del Barrio, R. Hermida, and N. Bagheradeh, "Low-power implementation of Mitchell's approximate logarithmic multiplication for convolutional neural networks," in *Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 2018, pp. 617–622.
- [60] M. S. Kim, A. A. Del Barrio, L. T. Oliveira, R. Hermida, and N. Bagherzadeh, "Efficient Mitchell's approximate log multipliers for convolutional neural networks," *IEEE Transactions on Computers*, vol. 68, no. 5, pp. 660–675, 2018.
- [61] S. S. Sarwar, S. Venkataramani, A. Raghunathan, and K. Roy, "Multiplier-less artificial neurons exploiting error resiliency for energyefficient neural computing," in *Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 2016, pp. 145–150.
- [62] S. Hashemi, R. I. Bahar, and S. Reda, "DRUM: A dynamic range unbiased multiplier for approximate applications," in *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2015, pp. 418–425.
- [63] U. Lotrič and P. Bulić, "Logarithmic multiplier in hardware implementation of neural networks," in *International Conference on Adaptive and Natural Computing Algorithms*. Springer, 2011, pp. 158–168.
- [64] H. Kim, M. S. Kim, A. A. Del Barrio, and N. Bagherzadeh, "A costefficient iterative truncated logarithmic multiplication for convolutional neural networks," in 2019 IEEE 26th Symposium on Computer Arithmetric (ARTH). IEEE, 2019, pp. 108–111.
- [65] M. S. Ansari, V. Mrazek, B. F. Cockburn, L. Sekanina, Z. Vasicek, and J. Han, "Improving the accuracy and hardware efficiency of neural networks using approximate multipliers," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 28, no. 2, pp. 317–328, 2019.
- [66] M. S. Ansari, B. F. Cockburn, and J. Han, "An improved logarithmic multiplier for energy-efficient neural computing," *IEEE Transactions* on Computers, vol. 70, no. 4, pp. 614–625, 2020.
- [67] J. Johnson, "Rethinking floating point for deep learning," arXiv preprint arXiv:1811.01721, 2018.
- [68] T.-B. Juang, C.-Y. Lin, and G.-Z. Lin, "Area-delay product efficient design for convolutional neural network circuits using logarithmic number systems," in *International SoC Design Conference (ISOCC)*. IEEE, 2018, pp. 170–171.
- [69] T.-B. Juang, H.-L. Kuo, and K.-S. Jan, "Lower-error and area-efficient antilogarithmic converters with bit-correction schemes," *Journal of the Chinese Institute of Engineers*, vol. 39, no. 1, pp. 57–63, 2016.
- [70] T.-B. Juang, P. K. Meher, and K.-S. Jan, "High-performance logarithmic converters using novel two-region bit-level manipulation schemes,"

- in Proceedings of 2011-1000 and a second second and Test: IEEE, 2011, pp. 1–4.
 17. Zhao, S. Dai, Ne Wenkatsan, M.-Yu, Liu, B. Khailany, B. Dally, and A. Anandkumar, "Low-precision training in logarithmic number system using multiplicative weight update," arXiv preprint arXiv:2106.13914, 2021.
- [72] S. C. ed, M. Liang, A. Guntoro, W. Stechele, and G. Ascheir "Efficient hardwares and a second star expresentation with arbitrary log-base," in *Proceedings of the International* Conference on Computer-Nided Design, 2018, pp. 1–8.
- [73] T.-Y. Lu, H.-H. Chin, H.-I. Wu, and R.-S. Tsay, "A very compact embedded CNN processor design based on logarithmic computing," arXiv preprint arXiv:2010.11686, 2020.
- [74] E. H. Lee, D. Miyashita, E. Chai, B. Murmann, and S. S. Wong, "LogNet: Energy-efficient neural networks using logarithmic computation," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASP)*. IEEE, 2017, pp. 5900–5904.
- [75] J. Xu, Y. Huan, L.-R. Zheng, and Z. Zou, "A low-power arithmetic element for multi-base logarithmic computation on deep neural networks," in *IEEE International System-on-Chip Conference (SOCC)*. IEEE, 2018, pp. 43–48.
- [76] J. Xu, Y. Huan, Y. Jin, H. Chu, L.-R. Zheng, and Z. Zou, "Basereconfigurable segmented logarithmic quantization and hardware design for deep neural networks," *Journal of Signal Processing Systems*, vol. 92, no. 11, pp. 1263–1276, 2020.
- [77] T. Ueki, K. Iwai, T. Matsubara, and T. Kurokawa, "Learning accelerator of deep neural networks with logarithmic quantization," in 2018 7th International Congress on Advanced Applied Informatics (IIAI-AAI). IEEE, 2018, pp. 634–638.

Number Systems for Deep Neural Network Architectures: A Survey

Ghada Alsuhli¹, Vasileios Sakellariou¹, Hani Saleh, Senior Member, IEEE,¹, Mahmoud Al-Qutayri, Senior Member, IEEE,¹, Baker Mohammad, Senior Member, IEEE,¹, and Thanos Stouraitis¹

33, 2011

928-931

- [47] M. G. Arnold, T. A. Bailey, J. J. Cupal, and M. D. Winkel, "On the cost effectiveness of logarithmic arithmetic for backpropagation training on simd processors," in Proceedings of International Conference on Neural Networks (ICNN'97), vol. 2. IEEE, 1997, pp. 933-936.
- [48] B. Parhami, "Computing with Electrical Engineering, vol.
- [49] D. Miyashita, E. H. Lee, ral networks using logarith arXiv:1603.01025, 2016.
- [50] A. Sanyal, P. A. Beerel, and with approximate logarithm Conference on Acoustics, IEEE, 2020, pp. 3122-3126.
- [51] S. A. Alam, J. Garland, an number systems: Beyond ba and Code Optimization (TAC
- [52] I. Kouretas and V. Palioura learning," in International C Technologies (MOCAST).
- [53] H. Saadat, H. Bokhari, an multipliers for approximate IEEE Transactions on Com and Systems, vol. 37, no. 11
- [54] J. N. Mitchell, "Computer logarithms," IRE Transaction 517, 1962.

Implementation methods at Jiawei Zhao, Steve Dai, Rangharajan Venkatesan, Brian Zimmer, Mustafa Ali, Ming-Yu Liu, Brucek Khailany, William J. Dally, Anima Anandkumar, 'LNS-Madam: Low-Precision Training in Logarithmic Number System using Multiplicative Weight Update" arXiv preprint 2106.13914v3

[55] Z. Babić, A. Avramović, and P. Bulić, "An iterative logarithmic

[56] M. S. Ansari, B. F. Cockburn, and J. Han, "A hardware-efficient

[57] R. Pilipović and P. Bulić, "On the design of logarithmic multiplier using

multiplier," Microprocessors and Microsystems, vol. 35, no. 1, pp. 23-

logarithmic multiplier with improved accuracy," in Design, Automation

& Test in Europe Conference & Exhibition (DATE). IEEE, 2019, pp.

Appeared in IEEE Trans. Comput. vol. 71, no. 12, pp. 3179-3190, 1 Dec. 2022, doi: 10.1109/TC.2022.3202747.

Authors are from NVIDIA, which has several related patent applications

Neural Network Accelerator Using Logarithmic-based Arithmetic Dally; William James ; et al.

uspto.report > / patents > / NVIDIA Corporation > / Patent 16/549683

Inference Accelerator Using Logarithmic-based Arithmetic Dally; William James ; et al.

and Test. IEEE, 2011, pp. 1-4.

J. Zhao, S. Dai, R. Venkatesan, M.-Y. Liu, B. Khailany, B. Dally, and

A. Anandkumar, "Low-precision training in logarithmic number system

using multiplicative weight update," arXiv preprint arXiv:2106.13914,

Liang, A. Guntoro, W. Stechele,

uspto.report > / patents > / NVIDIA Corporation > / Patent 16/750823





517, 1962.

LNS Addition via Gaussian Logs

Given
$$p = \log_b |P| < q = \log_b |Q|$$
:
1. Let $x = q - p$

Why it works: 1. x = log_b|Q/P|



LNS Addition via Gaussian Logs

Given
$$p = \log_b |P| < q = \log_b |Q|$$
:
1. Let $x = q - p$
2. Lookup $\Phi^+(x) = \log_b(1+b^x)$
or $\Phi^-(x) = \log_b |1-b^x|$

Why it works: 1. $x = \log_{b}|Q/P|$ 2. $\Phi^{+}(x) = \log_{b}(1+|Q/P|)$ $\Phi^{-}(x) = \log_{b}|1-|Q/P||$



LNS Addition via Gaussian Logs

Given
$$p = \log_b |P| < q = \log_b |Q|$$
:
1. Let $x = q - p$
2. Lookup $\Phi^+(x) = \log_b (1+b^x)$
or $\Phi^-(x) = \log_b |1-b^x|$
3. $t = p + \Phi(x)$

Why it works: 1. $x = \log_{b}|Q/P|$ 2. $\Phi^{+}(x) = \log_{b}(1+|Q/P|)$ $\Phi^{-}(x) = \log_{b}|1-|Q/P||$ 3. $t = \log_{b}(|P|(1\pm|Q|/|P|))$

Hardware:

- 1 min/max unit (so p<q)
- **1** subtractor
- 1 function approximation unit
- 1 adder



Gaussian Logs

Use: $\Phi^+(x)$ when signs are same, also known as $s_b(x) = \log_b(1+b^z)$ $\Phi^-(x)$ when signs are different, also known as $d_b(x) = \log_b|1-b^z|$ $\Phi(x)$ for generic Gaussian Log



Gaussian Logs

Use: $\Phi^+(x)$ when signs are same, also known as $s_b(x) = \log_b(1+b^z)$ $\Phi^-(x)$ when signs are different, also known as $d_b(x) = \log_b|1-b^z|$ $\Phi(x)$ for generic Gaussian Log

$\Phi^{-}(\mathbf{x})$ harder to approximate



Gaussian Logs

Use: $\Phi^+(x)$ when signs are same, also known as $s_b(x) = \log_b(1+b^z)$ $\Phi^-(x)$ when signs are different, also known as $d_b(x) = \log_b[1-b^z]$ $\Phi(x)$ for generic Gaussian Log

$\Phi^{-}(\mathbf{x})$ harder to approximate



Commutativity: only use x < 0 Φ(x) = Φ(-x) + x 1+X = X+1 = (1+1/X) X

Left Linear Taylor Interpolation



Right Linear Taylor Interpolation



Right Linear Taylor Interpolation



1. Partition range of x with non-uniform Δ [Lewis 1990]

F (precision)	15	17	19	21	23
Ф⁺ bits	0.3K	1K	2K	5K	10K
Φ ⁻ bits	1K	4K	10K	24K	60K
Problem: Φ ⁻ ta	kes mo	ost o	f the R	OM	

2. Separate non-interpolative computation of 1+b^x and log_b [Chen 2003] Problem: Number of stages proportional to precision

1. Partition range of x with non-uniform Δ [Lewis 1990]

F (precision)	15	17	19	21	23
Φ⁺ bits	0.3K	1K	2K	5K	10K
Φ ⁻ bits	1K	4K	10K	24K	60K
Problem: Φ ⁻ ta	kes mo	ost o	f the R	OM	

- 2. Separate non-interpolative computation of 1+b^x and log_b [Chen 2003] Problem: Number of stages proportional to precision
- 3. Use co-transformation to

a) Convert
$$\Phi^-$$
 with $\mathbf{x} = \mathbf{r}_b + \mathbf{r}_a$ to Φ^+ [Arnold 1998] :
 $\Phi^-(\mathbf{r}_b + \mathbf{r}_a) = \Phi^-(\mathbf{r}_b) + \Phi^+(\mathbf{r}_b + \Phi^-(\mathbf{r}_a) - \Phi^-(\mathbf{r}_b))$, where $\mathbf{r}_a > 0$ and $\mathbf{r}_b > 0$

1. Partition range of x with non-uniform Δ [Lewis 1990]

F (precision)	15	17	19	21	23
0[∓] bits	0.3K	1K	2K	5K	10K
Φ ⁻ bits	1K	4K	10K	24K	60K
Problem: Φ ⁻ ta	kes mo	ost o	f the R	OM	

- 2. Separate non-interpolative computation of 1+b^x and log_b [Chen 2003] Problem: Number of stages proportional to precision
- 3. Use co-transformation to

a) Convert Φ^- with $x = r_b + r_a$ to Φ^+ [Arnold 1998] : $\Phi^-(r_b + r_a) = \Phi^-(r_b) + \Phi^+(r_b + \Phi^-(r_a) - \Phi^-(r_b))$, where $r_a > 0$ and $r_b > 0$ b) Convert Φ^- with $x = r_b - r_a$ near 0 to Φ^- with arg away 0 [Coleman 2000]: $\Phi^-(r_b + r_a) = \Phi^-(r_b) + \Phi^-(r_b + \Phi^-(r_a) - \Phi^-(r_b))$, where $r_a < 0$ and $r_b > 0$

1. Partition range of x with non-uniform Δ [Lewis 1990]

F (precision)	15	17	19	21	23
Φ⁺ bits	0.3K	1K	2K	5K	10K
Φ ⁻ bits	1K	4K	10K	24K	60K
Problem: Φ ⁻ ta	kes mo	st o	f the R	OM	

- 2. Separate non-interpolative computation of 1+b^x and log_b [Chen 2003] Problem: Number of stages proportional to precision
- 3. Use co-transformation to

a) Convert
$$\Phi^-$$
 with $\mathbf{x} = \mathbf{r}_b + \mathbf{r}_a$ to Φ^+ [Arnold 1998] :
 $\Phi^-(\mathbf{r}_b + \mathbf{r}_a) = \Phi^-(\mathbf{r}_b) + \Phi^+(\mathbf{r}_b + \Phi^-(\mathbf{r}_a) - \Phi^-(\mathbf{r}_b))$, where $\mathbf{r}_a > 0$ and $\mathbf{r}_b > 0$
b) Convert Φ^- with $\mathbf{x} = \mathbf{r}_b - \mathbf{r}_a$ near 0 to Φ^- with arg away 0 [Coleman 2000]:
 $\Phi^-(\mathbf{r}_b + \mathbf{r}_a) = \Phi^-(\mathbf{r}_b) + \Phi^-(\mathbf{r}_b + \Phi^-(\mathbf{r}_a) - \Phi^-(\mathbf{r}_b))$, where $\mathbf{r}_a < 0$ and $\mathbf{r}_b > 0$

1. Partition range of x with non-uniform Δ [Lewis 1990]

F (precision)	15	17	19	21	23
Φ⁺ bits	0.3K	1K	2K	5K	10K
Φ ⁻ bits	1K	4K	10K	24K	60K
Problem: Φ ⁻ ta	kes mo	ost o	f the R	OM	

- 2. Separate non-interpolative computation of 1+b^x and log_b [Chen 2003] Problem: Number of stages proportional to precision
- 3. Use co-transformation to

a) Convert Φ^- with $x = r_b + r_a$ to Φ^+ [Arnold 1998]: $\Phi^-(r_b + r_a) = \Phi^-(r_b) + \Phi^+(r_b + \Phi^-(r_a) - \Phi^-(r_b))$, where $r_a > 0$ and $r_b > 0$ b) Convert Φ^- with $x = r_b - r_a$ near 0 to Φ^- with arg away 0 [Coleman 2000]: $\Phi^-(r_b + r_a) = \Phi^-(r_b) + \Phi^-(r_b + \Phi^-(r_a) - \Phi^-(r_b))$, where $r_a < 0$ and $r_b > 0$ $\Phi^-(r_a) = -r_a + \Phi(r_a) = \Phi^-(r_b) + \Phi^-(r_b - r_a + \Phi^-(r_a) - \Phi^-(r_b))$ $= \Phi^-(r_b) + \Phi^-(x_b + \Phi^-(r_a) - \Phi^-(r_b))$

Problem: still needs moderate-sized table of $\Phi^{-}(r_{a})$ and $\Phi^{-}(r_{b})$

1. Partition range of x with non-uniform Δ [Lewis 1989]

F (precision)	15	17	19	21	23
Φ⁺ bits	0.3K	1K	2K	5K	10K
Φ ⁻ bits	1K	4K	10K	24K	60K
Problem: Φ ⁻ ta	kes mo	ost o	f the R	OM	

- 2. Separate non-interpolative computation of 1+b^x and log_b [Chen 2003] Problem: Number of stages proportional to precision
- 3. Use co-transformation to

a) Convert Φ^{-} with $x = r_{b} + r_{a}$ to Φ^{+} [Arnold 1998]: $\Phi^{-}(r_{b} + r_{a}) = \Phi^{-}(r_{b}) + \Phi^{+}(r_{b} + \Phi^{-}(r_{a}) - \Phi^{-}(r_{b}))$, where $r_{a} > 0$ and $r_{b} > 0$ b) Convert Φ^{-} with $x = r_{b} - r_{a}$ near 0 to Φ^{-} with arg away 0 [Coleman 2000]: $\Phi^{-}(r_{b} + r_{a}) = \Phi^{-}(r_{b}) + \Phi^{-}(r_{b} + \Phi^{-}(r_{a}) - \Phi^{-}(r_{b}))$, where $r_{a} < 0$ and $r_{b} > 0$ $\Phi(-r_{a}) = -r_{a} + \Phi(r_{a}) = \Phi^{-}(r_{b}) + \Phi^{-}(r_{b} - r_{a} + \Phi^{-}(r_{a}) - \Phi^{-}(r_{b}))$ $= \Phi^{-}(r_{b}) + \Phi^{-}(x + \Phi^{-}(r_{a}) - \Phi^{-}(r_{b}))$ Problem: still needs moderate-sized table of $\Phi^{-}(r_{a})$ and $\Phi^{-}(r_{b})$ 4. Use higher-order co-transformation involving r_{a} , r_{b} , r_{c} ...: a) Recursive application of 3a) [Arnold 1997] or $\Phi^{-}(r_{c})$ Recursive application of 3b) [Coleman 2011] Problem: challenging to formalize

Formally verified error bounds for Φ^+ and Φ^- approximations



Lean4: A programming language and a proof assistant

• Built on dependent type theory

- Enables precise mathematical definitions and rigorous proof construction.
- Mathlib: A comprehensive mathematics library
 - Includes definitions and theorems crucial for formalizing diverse areas like calculus, algebra, and analysis.
- Metaprogramming capabilities
 - Enables custom tactics and automation to streamline the proof process.

Mathlib support for formalizing calculus-based proofs

• Limits of Functions

- Filter.Tendsto.add, Filter.Tendsto.rpow
- HasDerivAt.lhopital_zero_right_on_Ioo
- Derivatives
 - o deriv_add, deriv_sub, deriv_mul, deriv_div
- Continuity and Differentiability
 - Continuous.add, ContinuousOn.add, ContinuousAt.add
 - Differentiable.add, DifferentiableOn.add, DifferentiableAt.add
- Monotonicity and Antitonicity
 - o monotone_of_deriv_nonneg, antitone_of_deriv_nonpos
 - strictMono_of_deriv_pos, strictAnti_of_deriv_neg
- Tactics
 - fun_prop, continuity
 - o norm_num, positivity, linarith
 - o field_simp, ring_nf

First-order Taylor approximation



Error sources



Error of first-order Taylor approximation of Φ^+



• Consider the error function

$$E(\mathbf{i},\mathbf{r}) = \Phi(\mathbf{i}-\mathbf{r}) - (\Phi(\mathbf{i}) - \mathbf{r}\Phi'(\mathbf{i}))$$

as a function of two independent real-valued variable $i \leq 0$ and $r \geq 0$.

• Consider the error function

$$E(\mathbf{i},\mathbf{r}) = \Phi(\mathbf{i}-\mathbf{r}) - (\Phi(\mathbf{i}) - \mathbf{r}\Phi'(\mathbf{i}))$$

as a function of two independent real-valued variable $i \le 0$ and $r \ge 0$.

• Show that $E(\mathbf{i}, \mathbf{r})$ is a **strictly increasing** function of \mathbf{r} for fixed $\mathbf{i} \leq 0$.



• Consider the error function

$$E(\mathbf{i},\mathbf{r}) = \Phi(\mathbf{i}-\mathbf{r}) - (\Phi(\mathbf{i}) - \mathbf{r}\Phi'(\mathbf{i}))$$

as a function of two independent real-valued variable $i \le 0$ and $r \ge 0$.

- Show that $E(\mathbf{i}, \mathbf{r})$ is a **strictly increasing** function of \mathbf{r} for fixed $\mathbf{i} \leq 0$.
- Show that $E(\mathbf{i}, \mathbf{r})$ is a **strictly increasing** function of \mathbf{i} for fixed $0 \leq \mathbf{r}$.



• Consider the error function

$$E(\mathbf{i},\mathbf{r}) = \Phi(\mathbf{i}-\mathbf{r}) - (\Phi(\mathbf{i}) - \mathbf{r}\Phi'(\mathbf{i}))$$

as a function of two independent real-valued variable $i \le 0$ and $r \ge 0$.

- Show that $E(\mathbf{i}, \mathbf{r})$ is a **strictly increasing** function of \mathbf{r} for fixed $\mathbf{i} \leq 0$.
- Show that $E(\mathbf{i}, \mathbf{r})$ is a **strictly increasing** function of **i** for fixed $0 \leq \mathbf{r}$.
- The maximum error is reached when $\mathbf{i} = \mathbf{0}$ and $\mathbf{r} = \Delta$.



Error bound for first-order Taylor approximation of Φ^+

Lemma 1. For all
$$x \in (-\infty, 0]$$
,
 $|\Phi^+(x) - \hat{\Phi}^+_T(x)| \le E^+_{\Delta}(0) = \frac{\ln 2}{8}\Delta^2 + O(\Delta^4).$

lemma Ep_bound' (hi : i ≤ 0) (hr1 : 0 \leq r) (hr2 : r $\leq \Delta$) : |Ep i r| \leq Ep 0 Δ := by rw [abs_of_nonneg (Ep_r_nonneg hr1)] have ieq1 := Ep_i_monotoneOn hr1 hi Set.right_mem_Iic hi have ieq2 : Ep 0 r \leq Ep 0 Δ := Ep_r_monotoneOn hr1 (by linarith : 0 $\leq \Delta$) hr2 linarith

Error of first-order Taylor approximation of Φ^{-}

Y-Axis : Error of first-order Taylor approximation at x



Error of first-order Taylor approximation of Φ^{-}

Lemma 2. Suppose that $\frac{1}{\Delta} \in \mathbb{N}$ (e.g., $\Delta = 2^{-k}$ for some natural number k). Then for all $x \in (-\infty, -1]$,

$$|\Phi^{-}(x) - \hat{\Phi}_{T}^{-}(x)| \le -E_{\Delta}^{-}(-1) = (\ln 2)\Delta^{2} + O(\Delta^{3}).$$

lemma Em_bound' (hio : io < 0) (hi : i ≤ io) (hr1 : 0 ≤ r) (hr2 : r ≤ Δ) : |Em i r| ≤ Em io Δ := by have hi0 : i < 0 := by linarith rw [abs_of_nonneg (Em_r_nonneg hi0 hr1)] have ieq1 := Em_i_monotoneOn hr1 (by simp only [Set.mem_Iio]; linarith) (by simp only [Set.mem_Iio, hio]) hi have ieq2 : Em io r ≤ Em io Δ := Em_r_monotoneOn hio hr1 (by linarith : 0 ≤ Δ) hr2 linarith

Total error bound of first-order Taylor approximation

Theorem 1. Let ϵ be the machine-epsilon of the fixed-point representation of the LNS under consideration. Let $E_M^+ = E_{\Delta}^+(0)$, and $E_M^- = -E_{\Delta}^-(-1)$. Then

$$|\Phi(x) - \tilde{\Phi}_T(x)| < E_M + (2 + \Delta)\epsilon.$$

/- A model of fixed-point arithmetic -/
structure FixedPoint where
 ɛ : ℝ
 rnd : ℝ → ℝ
 hrnd : ∀ x, |x - rnd x| ≤ ε
 rnd_mono : Monotone rnd
 rnd_1 : rnd 1 = 1
 rnd_0 : rnd 0 = 0



Total error bound of first-order Taylor approximation

Theorem 1. Let ϵ be the machine-epsilon of the fixed-point representation of the LNS under consideration. Let $E_M^+ = E_{\Delta}^+(0)$, and $E_M^- = -E_{\Delta}^-(-1)$. Then

$$|\Phi(x) - \tilde{\Phi}_T(x)| < E_M + (2 + \Delta)\epsilon.$$

Subtraction requires special treatment for -1 < x < 0



Computing $\Phi^{-}(x)$ when x is close to 0 by interpolation is quite inaccurate

• Assume that x in (-1, 0). Fix a small value $\Delta_a > 0$. Write x = $r_b - r_a$ with $r_b = i \Delta_a$ for some integer i < 0 such that $r_b < x$ and $-\Delta_a \le r_a < 0$.

$$r_b = \left(\left\lceil \frac{x}{\Delta_a} \right\rceil - 1 \right) \Delta_a$$
$$r_a = r_b - x$$



- Assume that x in (-1, 0). Fix a small value $\Delta_a > 0$. Write x = r_b - r_a with r_b = i Δ_a for some integer i < 0 such that r_b < x and $-\Delta_a \le r_a < 0$. $r_b = \left(\left| \frac{x}{\Delta_a} \right| - 1 \right) \Delta_a$ $r_a = r_b - x$
- Compute Φ^{-} as $\Phi^{-}(x) = \Phi^{-}(r_{b} r_{a}) = \Phi^{-}(r_{b}) + \Phi^{-}(x + \Phi^{-}(r_{a}) \Phi^{-}(r_{b}))$.



• Assume that x in (-1, 0). Fix a small value $\Delta_a > 0$. $r_b = 0$. Write x = $r_b - r_a$ with $r_b = i \Delta_a$ for some integer i < 0 such that $r_b < x$ and $-\Delta_a \le r_a < 0$.

$$r_b = \left(\left\lceil \frac{x}{\Delta_a} \right\rceil - 1 \right) \Delta_a$$
$$r_a = r_b - x$$

- Compute Φ^- as $\Phi^-(x) = \Phi^-(r_b r_a) = \Phi^-(r_b) + \Phi^-(x + \Phi^-(r_a) \Phi^-(r_b))$.
- Requires 2 tables: T_a which contains all values of fixed-point values x in [- Δ_a , 0) and T_b which contains all integer multiples $r_b = i \Delta_a$ in (-1, Δ_a).



For x in (-1, 0) consider **2** cases:

- 1) x in $[-\Delta_a, 0]$. Take the value of $\Phi^-(x)$ directly from table T_a .
- 2) x in (-1, - Δ_a). Compute r_b and r_a , $\Phi^-(r_b)$ is taken from table T_b , $\Phi^-(r_a)$ is taken from table T_a , and $x + \Phi^-(r_a) \Phi^-(r_b) < -1$ is computed with other approximation techniques (e.g., Taylor approximation).



• Problem: Too many values in tables T_a and T_b when ε is small.

Example: $\varepsilon = 2^{-32}$, $\Delta_a = 2^{-16}$. Tables T_a and T_b have approximately 2^{16} values each.

- Problem: Too many values in tables T_a and T_b when ε is small.
 Example: ε = 2⁻³², Δ_a = 2⁻¹⁶.
 Tables T_a and T_b have approximately 2¹⁶ values each.
- Solution: Define $\Delta_b > \Delta_a$ and apply co-transformation technique twice. Compute r_c such that $x = r_c - r_{ab}$ with $r_c = j \Delta_b$ with j < 0 and $r_c < x$, $-\Delta_b \le r_{ab} < 0$. Apply co-transformation to r_{ab} and get $r_{ab} = r_b - r_a$, $r_b < r_{ab}$, $-\Delta_a \le r_a < 0$.

Requires one more table T_c . Tables T_a , T_b and T_c have approximately 2^{11} values each.

To compute $\Phi^{-}(x)$ consider **3** cases:

- 1) $-\Delta_a \le x < 0$: $\Phi^-(x)$ is taken from T_a .
- 2) x in $[-\Delta_b, -\Delta_a]$: $\Phi^-(x)$ is computed as $\Phi^-(r_b r_a) = \Phi^-(r_b) + \Phi^-(x + \Phi^-(r_a) \Phi^-(r_b))$.
- 3) x in $(-1, -\Delta_b)$: $\Phi^-(x)$ is computed as $\Phi^-(x) = \Phi^-(r_c - r_{ab}) = \Phi^-(r_c) + \Phi^-(x + \Phi^-(r_{ab}) - \Phi^-(r_c))$ Here $\Phi^-(r_c)$ is taken from table T_c and $\Phi^-(r_{ab}) = \Phi^-(r_b - r_a)$ is computed as in case 2.

To compute $\Phi^{-}(x)$ consider **3** 4 cases:

- 1) $-\Delta_a \le x < 0$: $\Phi^-(x)$ is taken from T_a .
- 2) x in $[-\Delta_b, -\Delta_a]$: $\Phi^-(x)$ is computed as $\Phi^-(r_b r_a) = \Phi^-(r_b) + \Phi^-(x + \Phi^-(r_a) \Phi^-(r_b))$.
- 3) x in $(-1, -\Delta_b)$: $\Phi^-(x)$ is computed as $\Phi^-(x) = \Phi^-(r_c - r_{ab}) = \Phi^-(r_c) + \Phi^-(x + \Phi^-(r_{ab}) - \Phi^-(r_c))$ Here $\Phi^-(r_c)$ is taken from table T_c and $\Phi^-(r_{ab}) = \Phi^-(r_b - r_a)$ is computed as in case 2.
- 4) x in (-1, $-\Delta_b$) and $r_{ab} < -\Delta_a$. In this case, we cannot compute $\Phi^-(r_{ab})$ as in case 2 because $r_{ab} + \Phi^-(r_a) \Phi^-(r_b)$ could be > -1. Take the value of $\Phi^-(r_{ab})$ directly from T_a .

Formalization revealed this additional case which was missing in our original informal proof.

Co-transformation error bound

Co-transformation formula with rounding:

$$\Phi^{-}(x) = \operatorname{rnd} \left(\Phi^{-}(r_{b}) \right) + \hat{\Phi}^{-} \left(x + \operatorname{rnd} \left(\Phi^{-}(r_{a}) \right) - \operatorname{rnd} \left(\Phi^{-}(r_{b}) \right) \right)$$

An approximation of $\Phi^{-}(x)$ for $x \leq -1$

A rounded value of $\Phi^{-}(x)$ A rounded value of $\Phi^{-}(x)$ stored in table T_b

A rounded value of $\Phi^{-}(x)$ stored in table T_b

Co-transformation error bound

Theorem 2. Let ϵ be the machine-epsilon of the fixed-point representation of the LNS under consideration and E_{Φ^-} be the error bound of interpolating of Φ^- in the range $(-\infty, 1]$. Assume also that $\Delta_a \ge 4\epsilon$ and $\Delta_b \ge 8\epsilon + 2E_{\Phi^-}$. The error bound of computing $\Phi^-(x)$ when $x \in (-1, 0)$ using the co-transformation technique is:

$$\Phi^{-}(-1 - E_{k_2}) - \Phi^{-}(-1) + E_{\Phi^{-}} + \epsilon$$

where

$$E_{k_2} = \Phi^-(-1 - 2\epsilon) - \Phi^-(-1) + E_{\Phi^-} + 2\epsilon.$$

theorem cotransformation_err_bound (fix : FixedPoint)(Φe : FunApprox Φm (Set.Iic (-1)))(ha : $0 < \Delta a$) (hb : $0 < \Delta b$)(h Δa : $4 * fix.\epsilon \le \Delta a$)(h Δa : $4 * fix.\epsilon \le \Delta a$)(h Δb : $8 * fix.\epsilon + 2 * \Phi e.err \le \Delta b$) : /- Δb should be large enough -/(h Δb : $8 * fix.\epsilon + 2 * \Phi e.err \le \Delta b$) : /- Δb should be large enough -/(h Δb : $2 * fix.\epsilon + \Phi m$ (-1 - 2 * fix. ϵ) - Φm (-1) + $\Phi e.err$ ($\Phi m x$ - Cotrans fix $\Phi e \Delta a \Delta b x$) \leq fix. $\epsilon + \Phi m$ (-1 - Ek2) - Φm (-1) + $\Phi e.err$

Co-transformation error bound (simplified)

Theorem 2. Let ϵ be the machine-epsilon of the fixed-point representation of the LNS under consideration and E_{Φ^-} be the error bound of interpolating of Φ^- in the range $(-\infty, 1]$. Assume also that $\Delta_a \ge 4\epsilon$ and $\Delta_b \ge 8\epsilon + 2E_{\Phi^-}$. The error bound of computing $\Phi^-(x)$ when $x \in (-1, 0)$ using the co-transformation technique is:

$$2E_{\Phi^-} + 5\epsilon$$

theorem cotransformation_err_bound' (fix : FixedPoint) $(\Phi e : FunApprox \Phi m (Set.Iic (-1)))$ $(ha : 0 < \Delta a)$ (hb : 0 < Δb) $(h\Delta a : 4 * fix.\epsilon \le \Delta a)$ $(h\Delta a : 4 * fix.\epsilon \le \Delta a)$ $(h\Delta b : 8 * fix.\epsilon + 2 * \Phi e.err \le \Delta b)$ $(h\Delta b : 8 * fix.\epsilon + 2 * \Phi e.err \le \Delta b)$ $(\Phi m x - Cotrans fix \Phi e \Delta a \Delta b x | \le 5 * fix.\epsilon + 2 * \Phi e.err := by$

Exhaustive verification



Exhaustive verification



Exhaustive verification



Results

- Formally verified error bounds for first order Taylor approximations of Φ⁺ and Φ⁻.
- Formally verified error bounds for co-transformation techniques for computing Φ⁻ near zero.
- A library of Lean 4 definitions and lemmas for error analysis of LNS: https://github.com/soarlab/RigorousErrorLNS

Future work

- Developing additional Lean 4 tactics and error-analysis lemmas.
- Tighter bounds for co-transformation.
- Error analysis for LNS variations:
 - Error Correction (EC)/quadratic
 - Other co-transformation [Arnold 1997]
 - unlike [Coleman 2000] does not use less accurate Φ⁻ interpolation
 - Bases other than two
 - \circ Guard bits
 - \circ Varying Δs
 - Table-less LNS design

Exhaustive verification with directed rounding



Metaprogramming - developing tactics

```
elab "get deriv" t:term loc:(deriv at)? : tactic => do
 let realType := Expr.const ``Real []
 withMainContext do
   let t ← Tactic.elabTerm t (some $ .forallE `x realType realType .default)
   let expr ← lambdaTelescope t (fun args e => toRExpr args e)
   let hyp := ← match loc with
                some loc => do
                 match loc with
                 (deriv at| at $x) => do
                     let x ← Tactic.elabTerm x (some realType)
                     mkAppM ``expr hasDerivAt #[expr, x]
                  (deriv at| within $s) => do
                     let s ← Tactic.elabTerm s (some realSetType)
                     mkAppM ``expr diff deriv on #[expr, s]
                  => throwUnsupportedSyntax
               => mkAppM ``expr diff deriv #[expr]
   let hypType ← inferType hyp
   liftMetaTactic fun mvarId => do
     let newId ← mvarId.assert `h concl (← mkAppM' hyp args)
     let (_, newId) ← newId.intro1P
     return args.toList.map Expr.mvarId! ++ [newId]
```