

Supplementary Information

Fig S1. Efficacy and specificity of IdeS vs papain proteases for Fc cleavage, trypsin vs chymotrypsin digestion for MS, and gel slicing for MS. (A) 10 µg of representative proteins were digested with the indicated amount of papain or 67 Units of IdeS for 1 hour at 37°C, then analyzed by SDS-PAGE. (B) The serum HCAb fraction from an immunized llama was affinity purified on GFP Sepharose, then digested with papain or IdeS, or left undigested. Remaining resin-bound V_HH or HCAb was eluted in LDS and separated by SDS-PAGE before trypsin digest and LC-MS/MS analysis. Llama-Magic was used to generate nanobody candidate lists for each sample, which were then searched for GFP nanobody sequences previously identified from this serum. The positions of identified nanobodies in the total candidate lists are graphed. (C) Trypsin and chymotrypsin digestions of purified V_HH bands were analyzed by MS, and data searched either separately or together for nanobody identification. For each of the three searches, a list of fully unique matching sequences was generated. This list was filtered so that each CDR3 group is represented by the member with the largest number of DNA sequencing counts. This filtered list was sorted from high to low MS coverage. The number of identified CDR3 groups is plotted against the MS coverage percentage of the corresponding CDR3 region or total V_HH sequence. (D) V_HH purified against SARS-CoV-2 spike was separated by SDS-PAGE and the band sliced into 5 pieces analyzed separately by MS, or treated as a single band (“Non-sliced”). The cumulative count of unique PSMs over the 5 slices (“Slice 1-5”) was significantly higher than the unsliced control. For two example V_HH sequences, relative intensities of three peptides covering CDR1, CDR2, and CDR3 were compared across the 5 slices. Differing enrichment patterns are seen consistently across the 3 peptides measured for each nanobody, assisting matching of CDR peptides to a single V_HH sequence.

Fig S2. Comparing V_HH coverage of nested PCR and single-step hinge PCR. (A) Schematic comparison of traditional nested PCR and single-step hinge PCR protocols. (B) CALL001/CALL002 PCR amplicons from llama or alpaca lymphocyte cDNA were sequenced by PacBio, and the most abundant V_HH N-term and C-term 18 bp sequences were identified to assess potential priming locations. Percentage of the total analyzed population is indicated for each 18bp sequence. (C) Hinge PCR primers were searched against lymphocyte mRNA sequences to estimate relative coverage. The top 10 complementary sequence hits are shown, with percentage of total matches. Exact primer matches are indicated in blue, and mismatched nucleotides labeled in red. SH-rev matches also included off-target non-IgG sequences, indicated by asterisks.

Figure S3. SPR sensorgrams for GFP nanobodies. Representative sensorgrams from SPR binding experiments of nanobody analytes flowed over an immobilized GFP ligand are shown. Curves are from parallel binding kinetics experiments (LaG94-2, LaG94-4, LaG94-6, LaG94-10, LaG94-12, and LaG94-18) on a ProteOn instrument, or single-cycle kinetics experiments on a Biacore 8K instrument (injected nanobody concentrations = 0.1, 0.3, 1, 3, and 10 nM). Curves from a Langmuir 1:1 binding model were fitted.

Figure S4. SPR sensorgrams for tdTomato nanobodies. Representative sensorgrams from SPR binding experiments of nanobody analytes flowed over an immobilized tdTomato ligand are shown. Curves are from single-cycle kinetics experiments on a Biacore 8K instrument (injected nanobody concentrations = 0.1, 0.3, 1, 3, and 10 nM). Curves from a Langmuir 1:1 binding model were fitted.

Figure S5. SPR sensorgrams for GST nanobodies. Representative sensorgrams from SPR binding experiments of nanobody analytes flowed over immobilized GST-TEV are shown. Curves are from single-cycle kinetics experiments on a Biacore 8K instrument (injected nanobody concentrations = 0.1, 0.3, 1, 3, and 10 nM). Curves from a Langmuir 1:1 binding model were fitted.

Figure S6. SPR sensorgrams for IgG nanobodies. Representative sensorgrams from SPR binding experiments of nanobody analytes flowed over immobilized goat, mouse, or rabbit IgG ligands are shown. Curves are from single-cycle kinetics experiments on a Biacore 8K instrument (injected nanobody concentrations = 0.1, 0.3, 1, 3, and 10 nM). Curve fits are from a Langmuir 1:1 binding model, or a two-step binding model (asterisks).

Figure S7. Screening GFP nanobodies by immunofluorescence. (A) Schematic of nanobody conjugation and staining procedure. Nanobodies were labeled with Alexa Fluor 568 and used to stain brain sections of Thy1-GFPM mice. (B) Fluorescent imaging was performed on GFP or nanobody-AF568 signal. (C) Quantitative comparisons of relative GFP and nanobody fluorescence. Fluorescence was normalized to maximum fluorescence in each channel. Plots were fitted with simple linear regression, and Pearson correlation coefficients (r) are indicated. Scale bars are 100 μm . Micrographs and quantification panels for LaG94-12, LaG94-15, and LaG94-18 are reproduced here from Fig. 4C to aid in comparison between samples.

Figure S8. Screening tdTomato nanobodies by immunofluorescence. (A) Schematic of nanobody conjugation and staining procedure. Nanobodies were labeled with Alexa Fluor 647 and used to stain sagittal brain sections of ChAT-Cre::Ai14 mice, with tdTomato expressed in cholinergic neurons. (B) Fluorescent imaging was performed on tdTomato or nanobody-AF647 signal. (C) Quantitative comparison of relative tdTomato and nanobody fluorescence. Fluorescence was normalized to maximum fluorescence in each channel. Plots were fitted with simple linear regression, and Pearson correlation coefficients (r) are indicated. Scale bars are 100 μm . Micrographs and quantification panels for LaTdT-1 and LaTdT-39 are reproduced here from Fig. 4D to aid in comparison between samples.

Figure S9. Complementation of GFP nanobodies by immunofluorescence (LaG2). (A) Schematic of nanobody complementation IF procedure. Nanobodies were labeled with Alexa Fluor 568 (LaG2) or 647 (Complement) and used to stain brain sections of Thy1-GFPM mice in sequence. (B) Fluorescent imaging was performed on GFP, nanobody-AF568 (first nanobody, LaG2), or nanobody-AF647 signal (second nanobody, indicated at left). Scale bars are 100 μm .

Figure S10. Complementation of GFP nanobodies by immunofluorescence (LaG16). Nanobodies were labeled with Alexa Fluor 568 (LaG16) or 647 (Complement) and used to stain brain sections of Thy1-GFPM mice in sequence. Fluorescent imaging was performed on GFP, nanobody-AF568 (first nanobody, LaG16), or nanobody-AF647 signal (second nanobody, indicated at left). Scale bars are 100 μm .

Figure S11. Complementation of GFP nanobodies by immunofluorescence (LaG41). Nanobodies were labeled with Alexa Fluor 568 (LaG41) or 647 (Complement) and used to stain brain sections of Thy1-GFPM mice in sequence. Fluorescent imaging was performed on GFP, nanobody-AF568 (first nanobody, LaG41), or nanobody-AF647 signal (second nanobody, indicated at left). Scale bars are 100 μm .

Figure S12. Complementation of GFP nanobodies by immunofluorescence (LaG94-10,14,15). Nanobodies were labeled with Alexa Fluor 568 or 647 and used to stain brain sections of Thy1-GFPM mice in sequence. Sections were first stained with (A) LaG94-10, (B) LaG94-14, or (C) LaG94-15. Fluorescent imaging was performed on GFP, nanobody-AF568 (first nanobody), or nanobody-AF647 (second nanobody, indicated at left) signal. Scale bars are 100 μm .

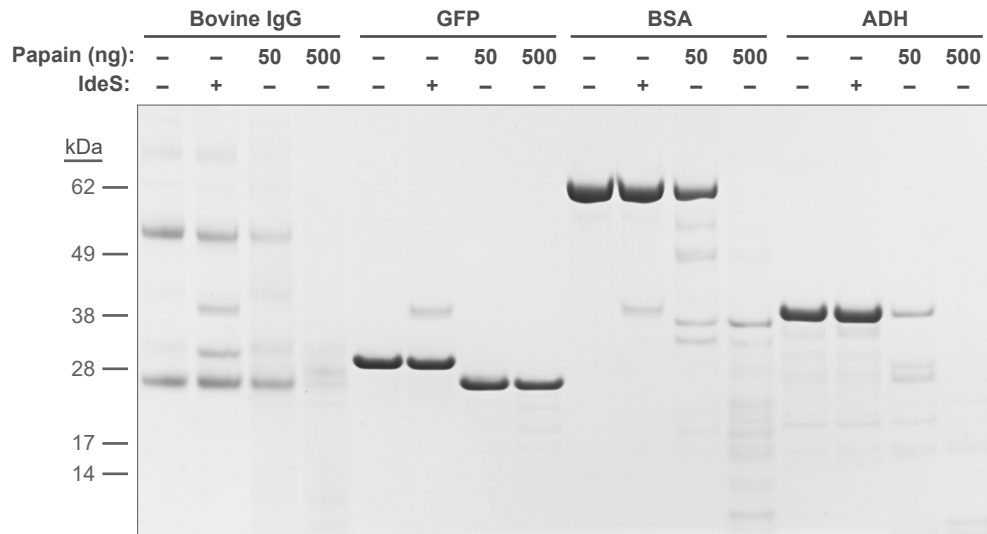
Figure S13. Screening GFP nanobody dimers by immunofluorescence. GFP nanobody dimers were labeled with Alexa Fluor 647 and used to stain brain sections of Thy1-GFPM mice. (A) Fluorescent imaging was performed on GFP or nanobody-AF647 signal. (B) Quantitative comparisons of relative GFP and nanobody fluorescence. Fluorescence was normalized to maximum fluorescence in each channel. Plots were fitted with simple linear regression, and Pearson correlation coefficients (r) are indicated. Scale bars are 100 μm . Micrographs and quantification panels for LaG 94-1--94-14 are reproduced here from Fig. 4E to aid in comparison between samples.

Figure S14. Screening IgG nanobodies by immunofluorescence. (A) Schematic of nanobody conjugation and staining procedure. IgG nanobodies were labeled with Alexa Fluor 568 or 647 and used to stain brain sections of Thy1-GFPM probed with an anti-GFP antibody of the corresponding species (mouse or rabbit). For goat IgG nanobodies, somatostatin tdTomato mouse brain sections were stained with goat anti-tdTomato. Fluorescent imaging was performed on GFP/tdTomato or nanobody signal using (B) mouse IgG nanobodies, (C) rabbit IgG nanobodies, or (D) goat IgG nanobodies. Traditional rabbit or goat secondaries were also used as controls. Quantitative comparisons of relative GFP or tdTomato and nanobody/secondary fluorescence are plotted at right. Fluorescence was normalized to maximum fluorescence in each channel. Plots were fitted with simple linear regression, and Pearson correlation coefficients (r) are indicated. Scale bars are 100 μm . Micrographs and quantification panels for LaMIgG-8, LaMIgG-14, and Goat Anti Mouse are reproduced from Fig. 5A in (B); LaRIgG-1 and LaRIgG-2 from Fig. 5B in (C); and LaGIgG-12 and Donkey Anti Goat from Fig. 5C in (D), to aid in comparison between samples.

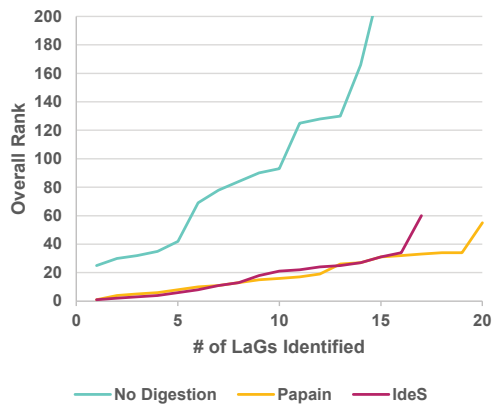
Figure S15. Species and isotype specificity of IgG nanobodies. Nanobodies were labeled with Alexa Fluor 488 and used to probe immobilized IgG of the indicated species, fragment, and/or isotype. An anti-GFP nanobody (LaG) was used as a negative control.

Figure S1

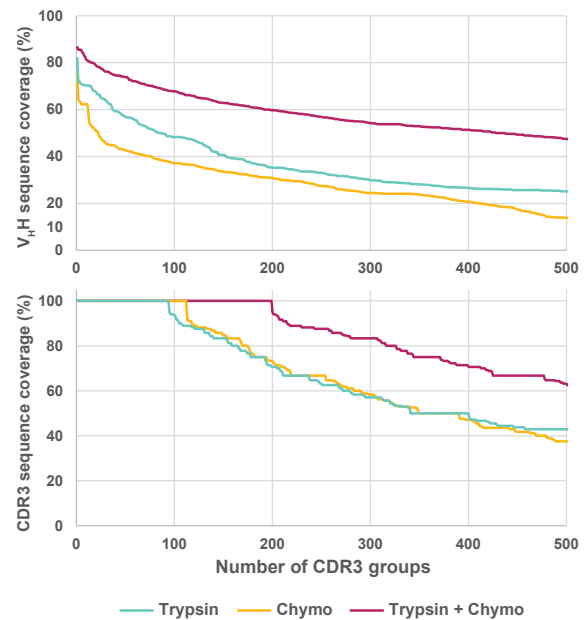
A



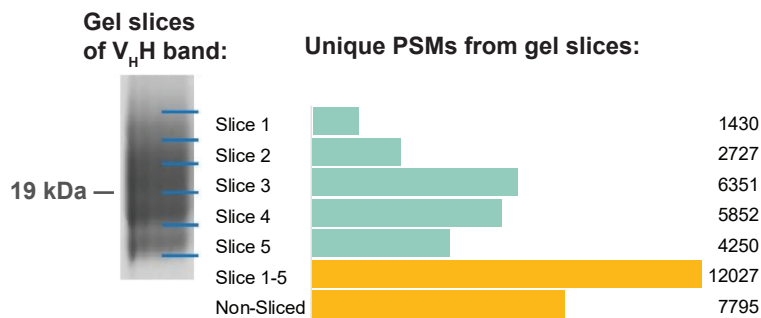
B



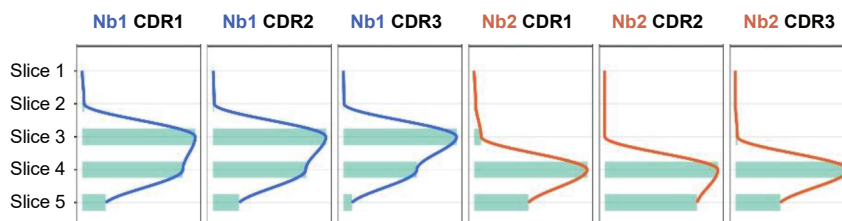
C



D

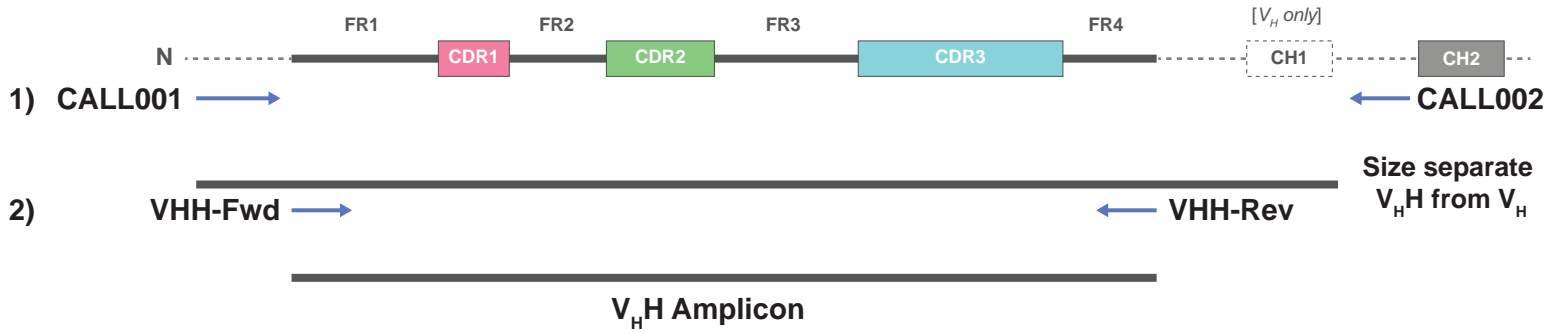


Relative intensity of CDR peptides from gel slices:

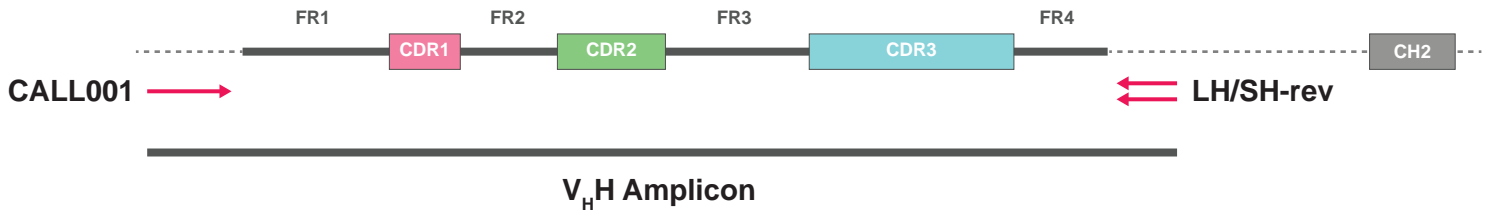


A Figure S2

Version 1.0 Nested PCR



Version 2.0 Hinge PCR



B

Llama V_HH

N-term V _H H (FR1)		C-term V _H H (FR4)	
caggtgcagctggtggag	8.1%	accagggtcaccgtctcc	62.5%
caggtgcagctggttagag	7.6%	accctgggtcaccgtctcc	6.6%
caggtgcagctggtggag	4.4%	accagggtcaccgtctct	2.9%
gaggtgcagctcgtggag	4.0%	actcaggtcaccgtctcc	2.6%
caggtgcaggtggtggaa	3.0%	atcctgggtcaccatctcc	1.9%
caggtgcagctcgtggag	2.5%	acccaagtcaccgtctcc	1.7%
caggtgcagctggcagag	2.4%	accagggtcaccgtcgcc	1.6%
caggtgcacctggttagag	2.3%	acgcaggtcaccgtctcc	1.2%
caggtgcaactggtggag	2.0%	atccaggtcaccgtctcc	1.1%
caggtgcagttggttagag	1.6%	accagggtcactgtctcc	1.1%
caggtgcatttggtggag	1.4%	accagggtcaccgtctcg	1.1%
caggtgcagctcgtggag	1.4%	accgggtcaccgtttcc	1.0%
caggtgcaactggttagag	1.3%	accagggttaccgtctcc	1.0%
caggtgcagctggttagaa	0.7%	accagggtcaccgtttcc	0.9%
gaggagccactcgtggag	0.7%	atccaggtcaccgtctcc	0.9%
gcggtgcagctcgtggag	0.7%	accgggtcaccgtctcc	0.7%
caggtgcaggtcgtggag	0.6%	accagggtcaccgtctca	0.7%
gaggtgcagctggtggag	0.6%	accctgggtcaccgtttcc	0.7%
caggtgcaggtggtggag	0.6%	accctgggtcactgtctcc	0.6%
caggtgcacctggtggag	0.6%	accagggtgaccgtctcc	0.6%
Top 20 Total:	47%	Top 20 Total:	92%

Alpaca V_HH

N-term V _H H (FR1)		C-term V _H H (FR4)	
caggtgcagctggtggag	12.1%	accagggtcaccgtctcc	63.7%
cagttgcagctcgtggag	5.8%	accctgggtcaccgtctcc	4.3%
caggtgcagctcgtggag	4.5%	accagggtcaccgtctct	3.2%
gaggtgcagctcgtggag	4.2%	actcaggtcaccgtctcc	3.2%
caggtgcagttggtggag	2.0%	accagggtcactgtctcc	1.7%
caggtgcaactggtggag	1.6%	accagggtcaccgtcaacc	1.7%
cagttgcaactcgtggag	1.5%	atccaggtcaccgtctcc	1.6%
caggtgcagctcgtggag	1.3%	accagggtcaccgtttcc	1.5%
gaaatcacagttcgtggag	1.2%	acccaagtcaccgtctcc	1.4%
gaggtgcagctggtggag	1.0%	accagggtcaccgtcgcc	1.3%
cagttgcagttcgtggag	0.8%	accgggtcaccgtctcc	1.3%
caggtgcacctggtggag	0.6%	accagggtcactgtctcc	1.2%
caggtgcagctggtggag	0.6%	accagggtcaccgtctcg	1.1%
caggtgcagctggtggag	0.6%	accagggttaccgtctcc	0.9%
gcggtgcagctcgtggag	0.6%	accctgggtcaccgtttcc	0.9%
gaggtgcagctcgtggag	0.5%	acgcaggtcaccgtctcc	0.7%
caggtgacctggtggag	0.5%	accctgggtcaccgtctct	0.6%
cagttgcaggtcgtggag	0.5%	accagggtgaccgtctcc	0.6%
caggtgcaactcgtggag	0.5%	accagggtcaccgtgtcc	0.5%
caggagcaagtgaaggag	0.5%	accagggtcaccgtctca	0.5%
Top 20 Total:	41%	Top 20 Total:	92%

C

CALL001 (N-term Leader):

gtcctggctgctcttttacaagg	25.1%
gtcctggctgctctactacaagg	23.3%
gtcctggctgctcttttacaagg	21.4%
gtcctggctgctcttttacaagg	5.3%
gtcctggctgctcttttacaagg	3.9%
gtcctggctgctcttttacaagg	3.2%
gtactggctgctctactacaagg	1.7%
gtgctggcggcccttctgctagg	1.2%
gtcctggctgctctcttacaagg	0.9%
gtcctggctgctctactacagg	0.8%
Top 10 Total:	87%

SH-rev (Short Hinge):

gcaccacagcgaagaccaccag	67.5%
*gcaccacagcgaagaccaccac	8.5%
*gccttacaggaggaccaccag	4.7%
gcaccacagcgaagaccaccag	3.2%
*gcaccacagcgaagaccaccag	2.8%
*tcaccacagtgacgaccaccag	1.3%
gcaccacagcgaagaccaccag	1.0%
*gcaccacagcgaagaccaccag	0.9%
*aaaccacagcgaagaccaccag	0.7%
*gcaccacagtgaaattccctag	0.6%
Top 10 Total:	91%

LH-rev (Long Hinge):

cccaagacacccaaaaccacaaccac	86.9%
cccaagacacccaaaaccacaaccac	1.4%
cccaagcaccacccaaaaccacaaccac	1.2%
cccaagacacccaaaaccacaaccac	0.8%
cccaagacacccaaaaccacaaccac	0.7%
cccaagacacccaaaaccacaaccac	0.7%
cccaagacacccaaaaccacaaccac	0.5%
cccaagacacccaaaaccacaaccac	0.5%
cccaagacacccaaaaccacaaccac	0.5%
cccaagacacccaaaaccacaaccac	0.5%
Top 10 Total:	94%

*Adjacent to non-FR4/hinge sequence

Figure S3

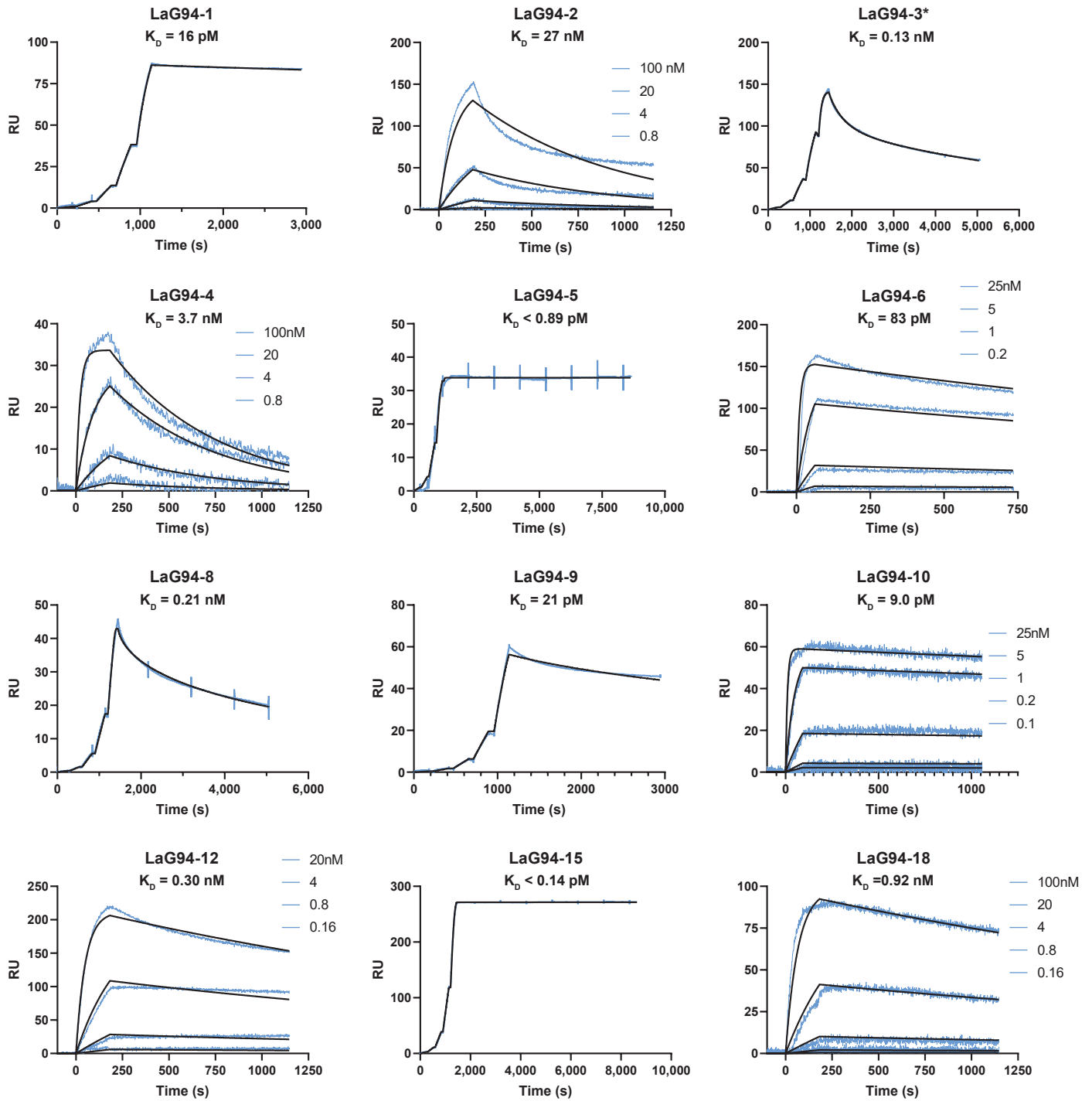


Figure S4

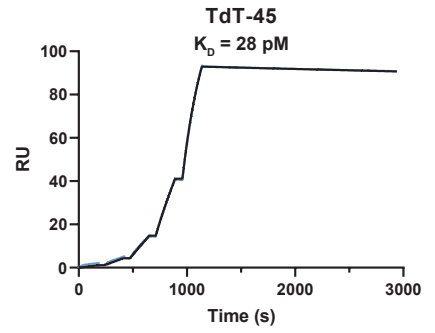
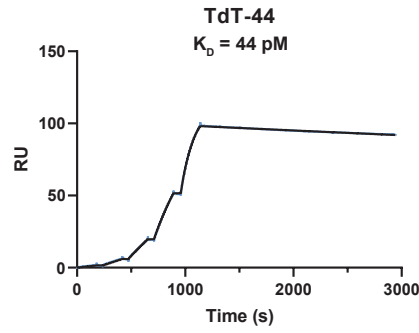
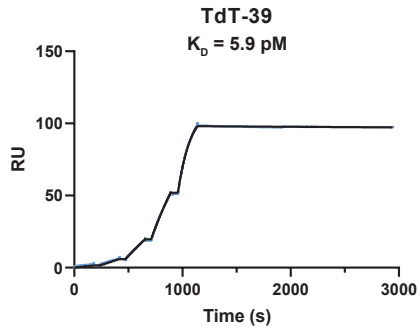
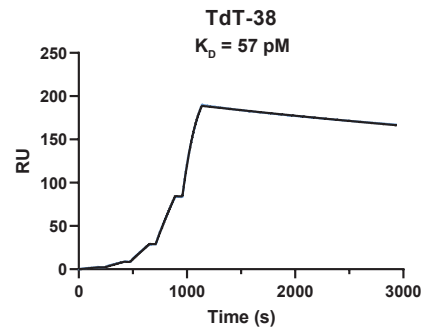
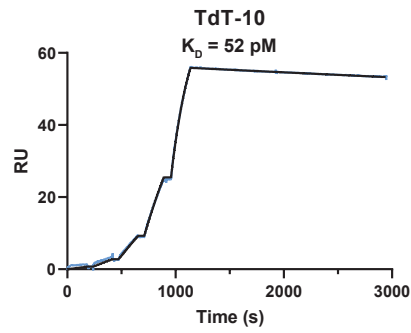
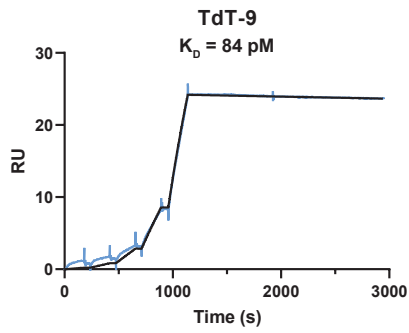
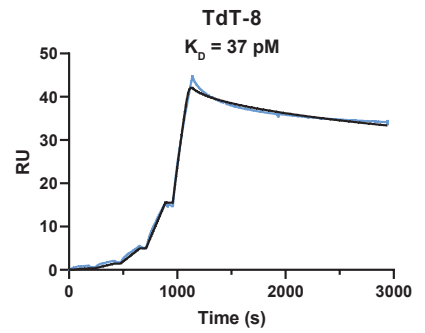
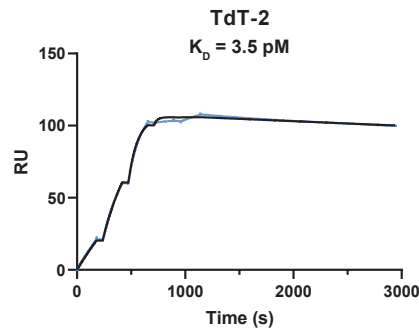
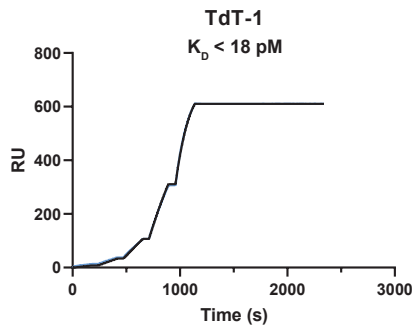


Figure S5

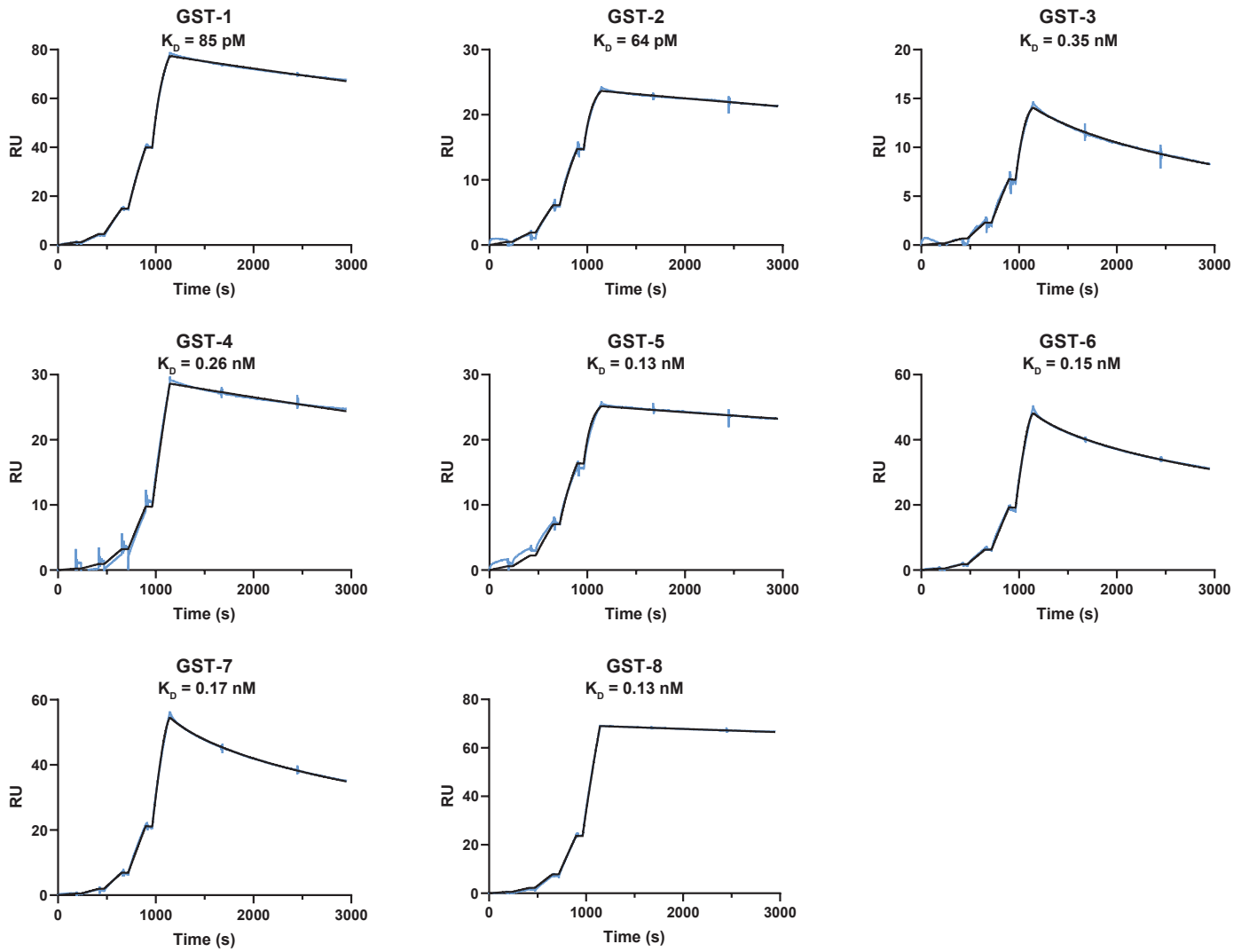


Figure S6

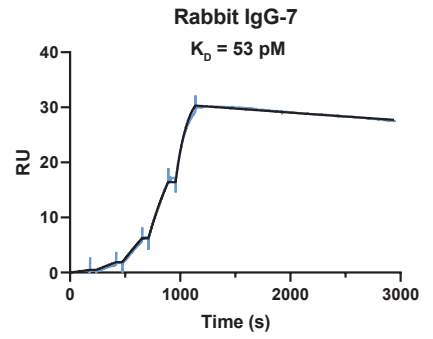
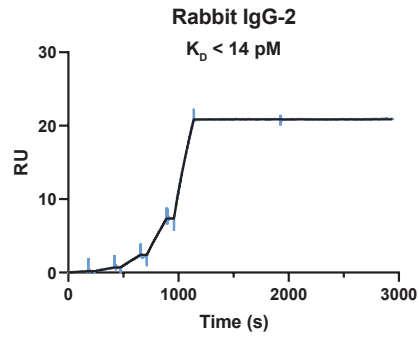
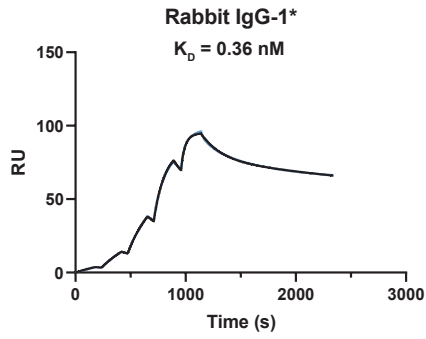
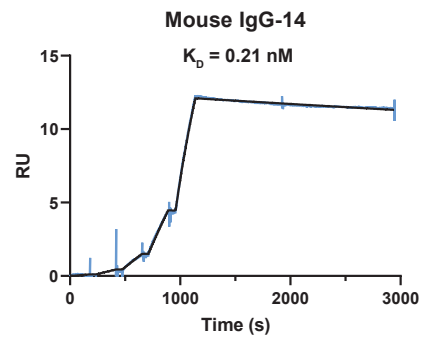
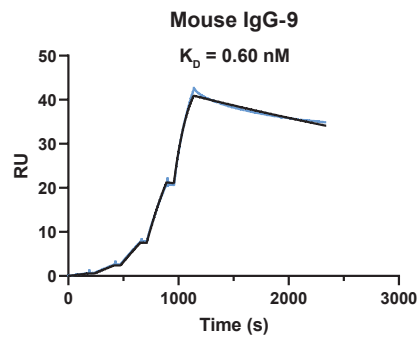
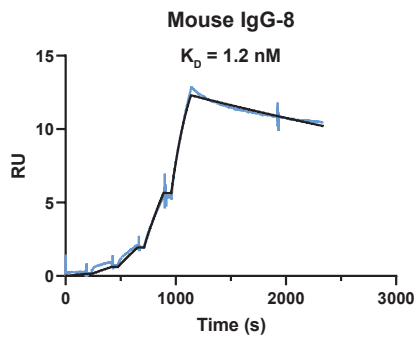
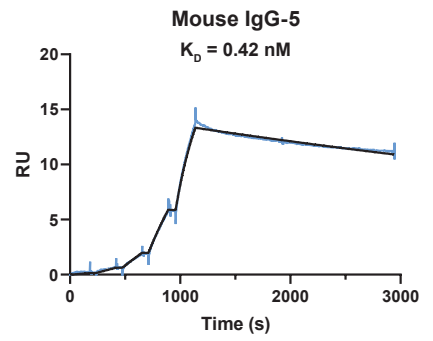
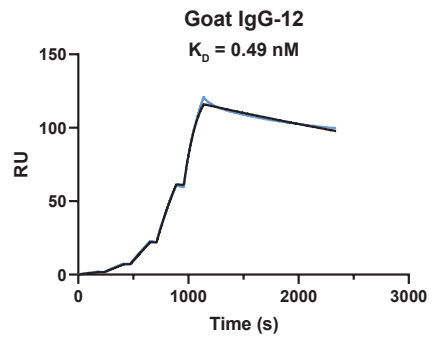
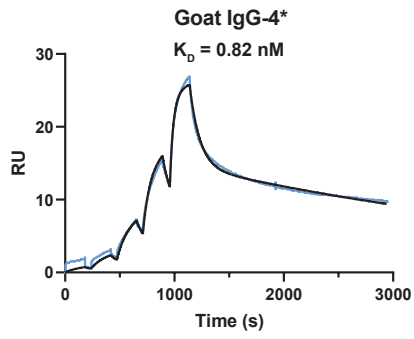


Figure S7

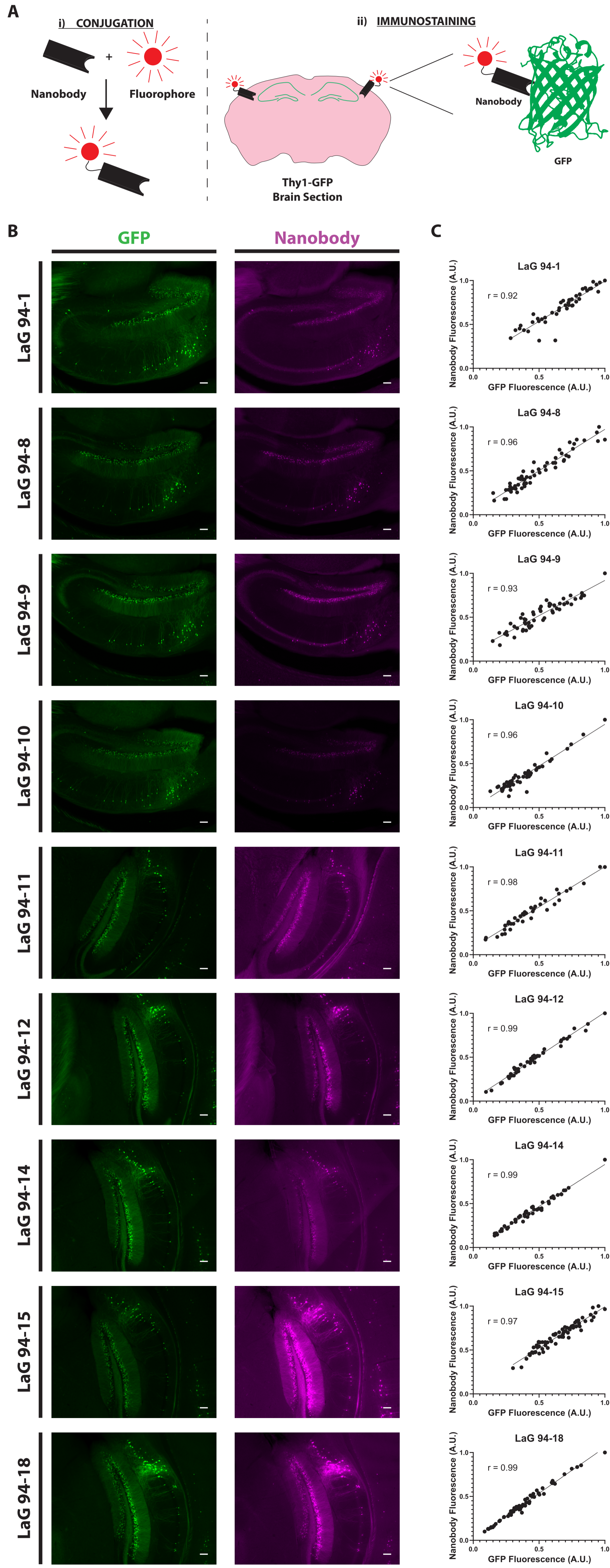


Figure S8

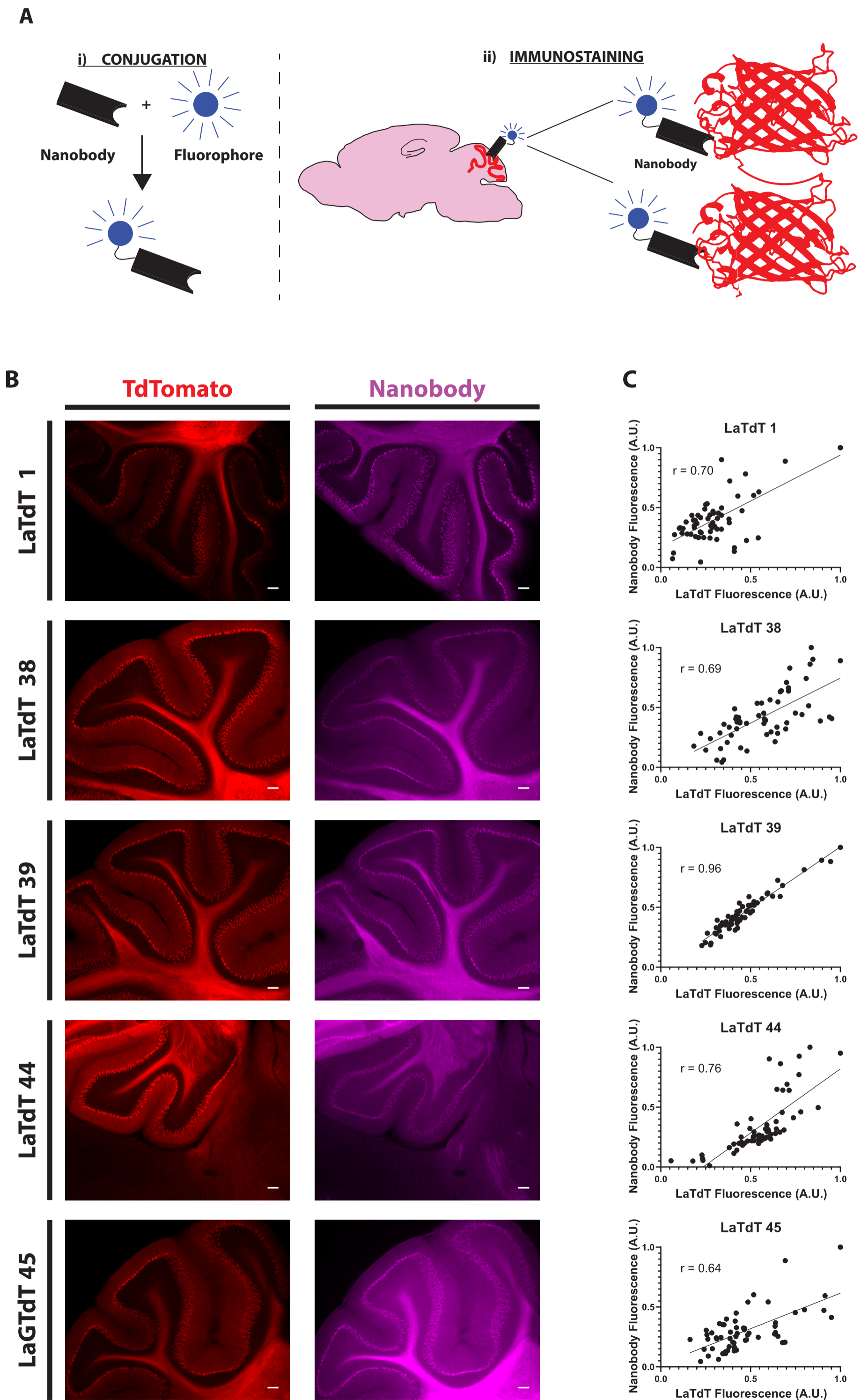


Figure S9

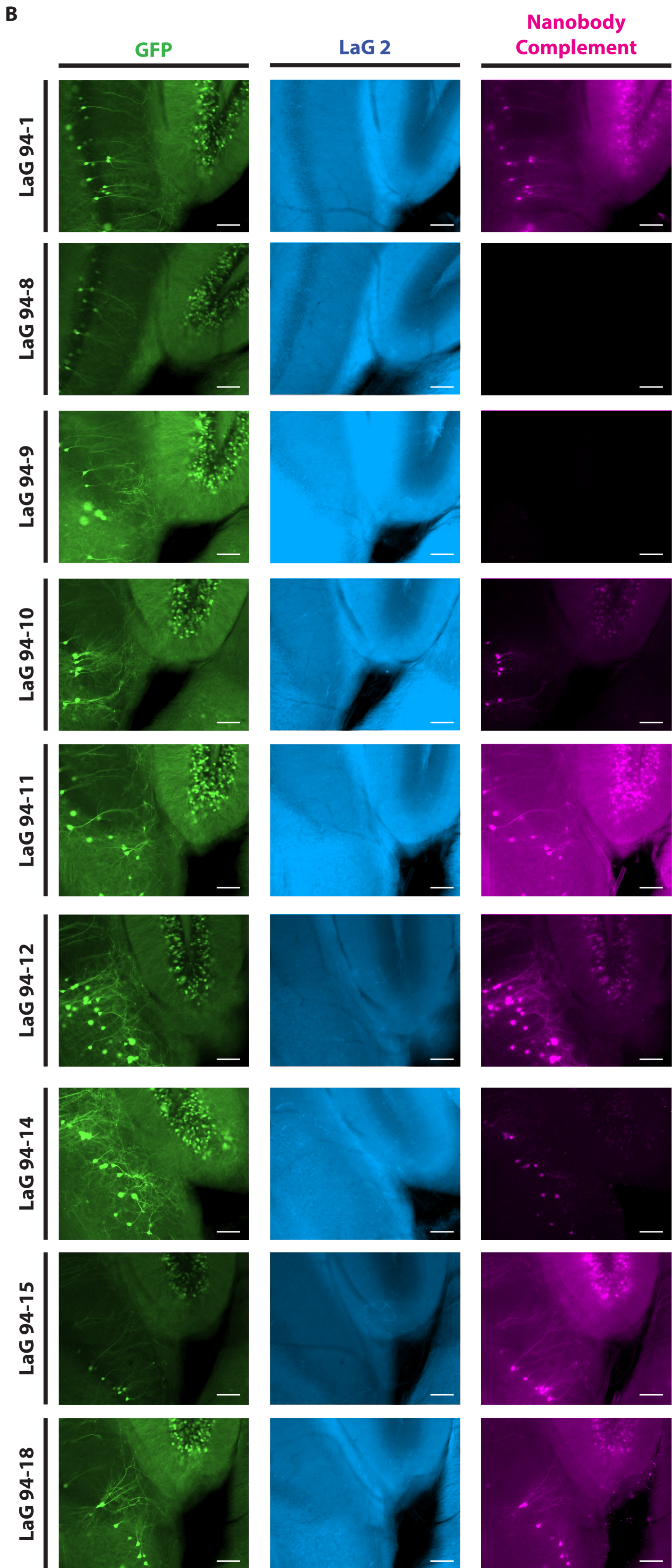
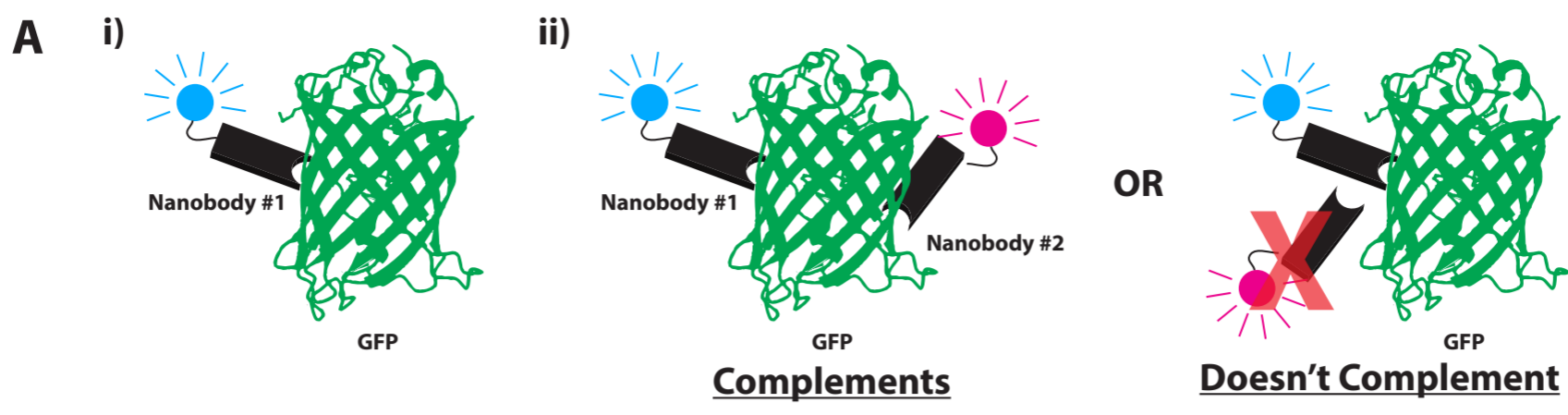


Figure S10

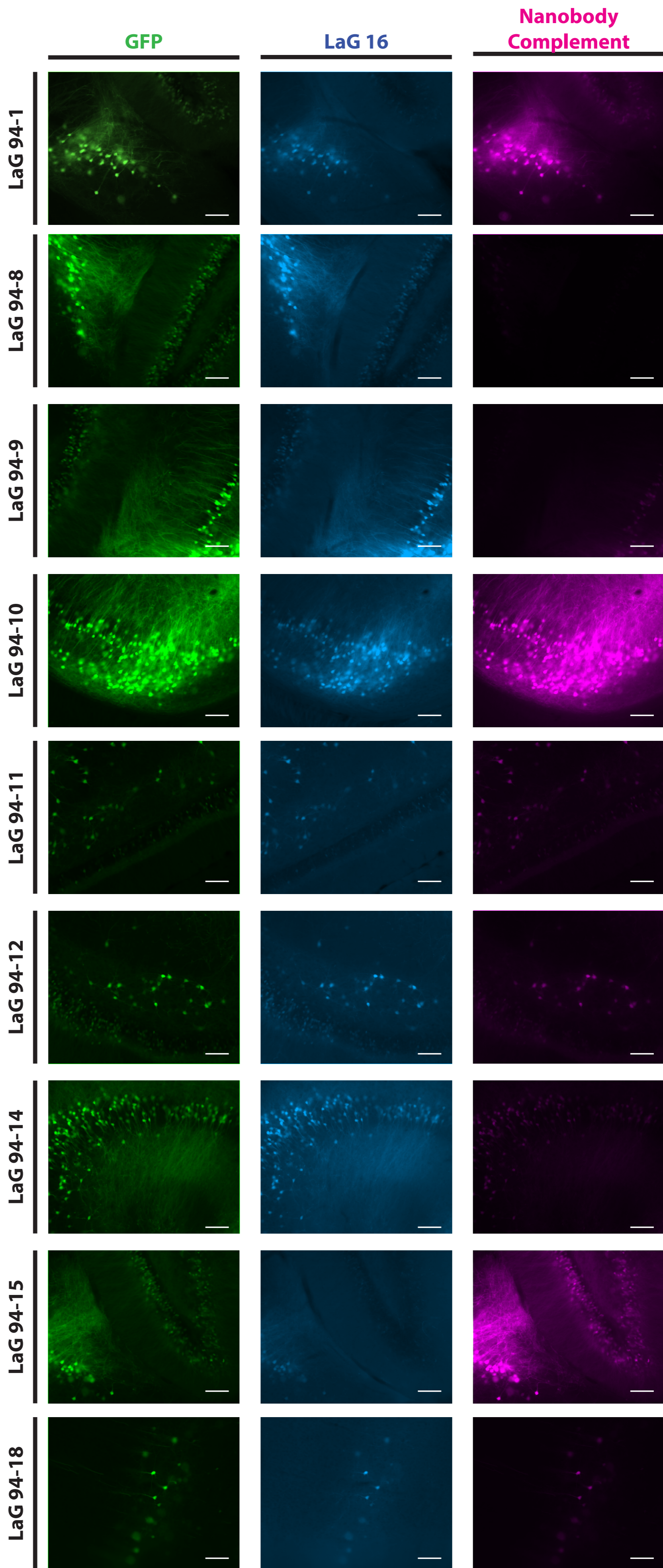


Figure S11

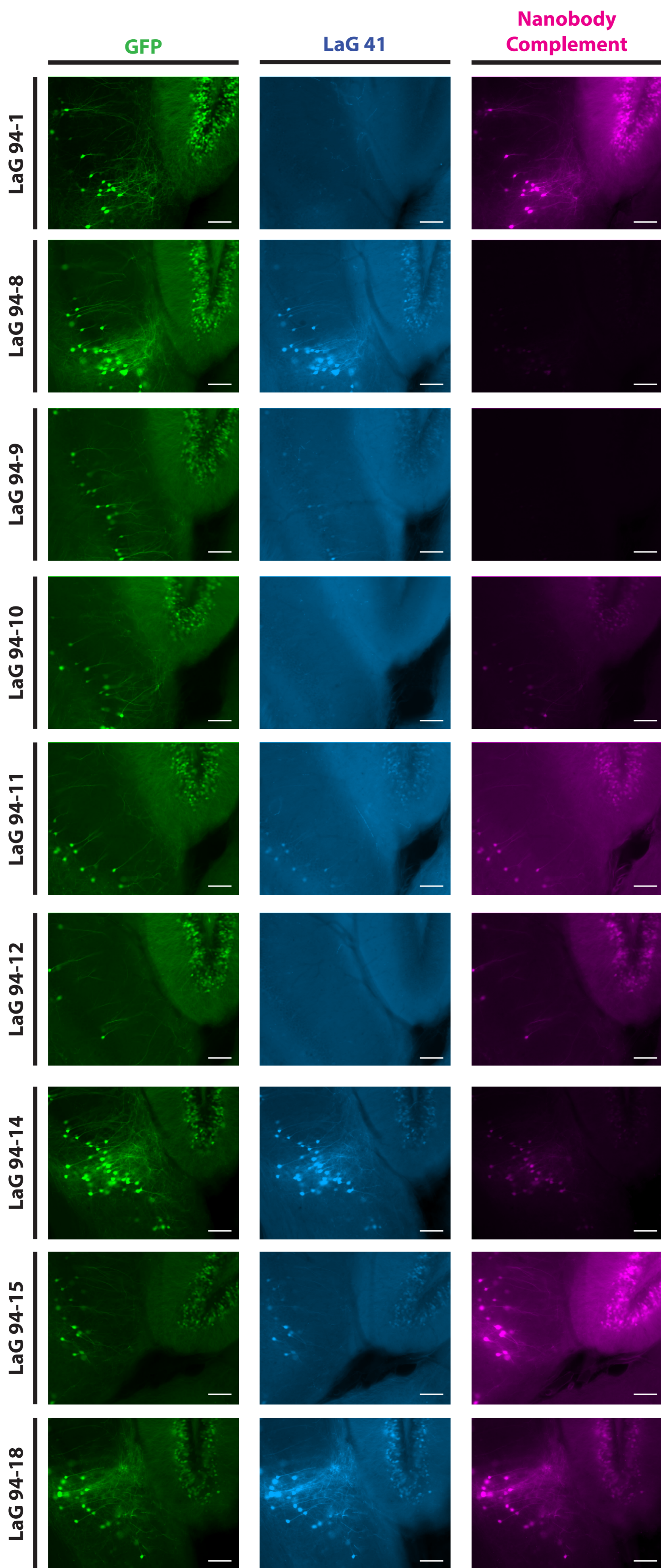
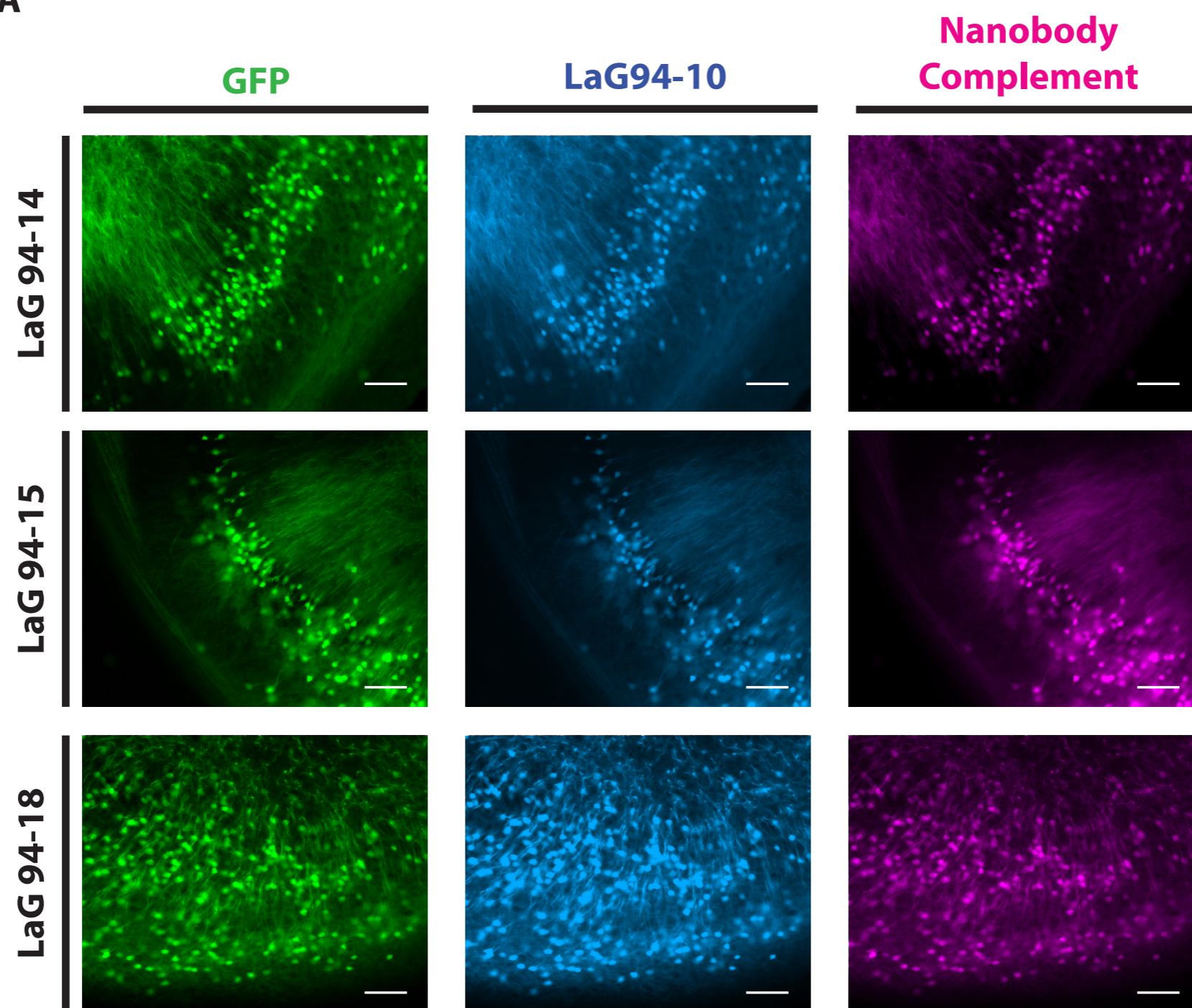
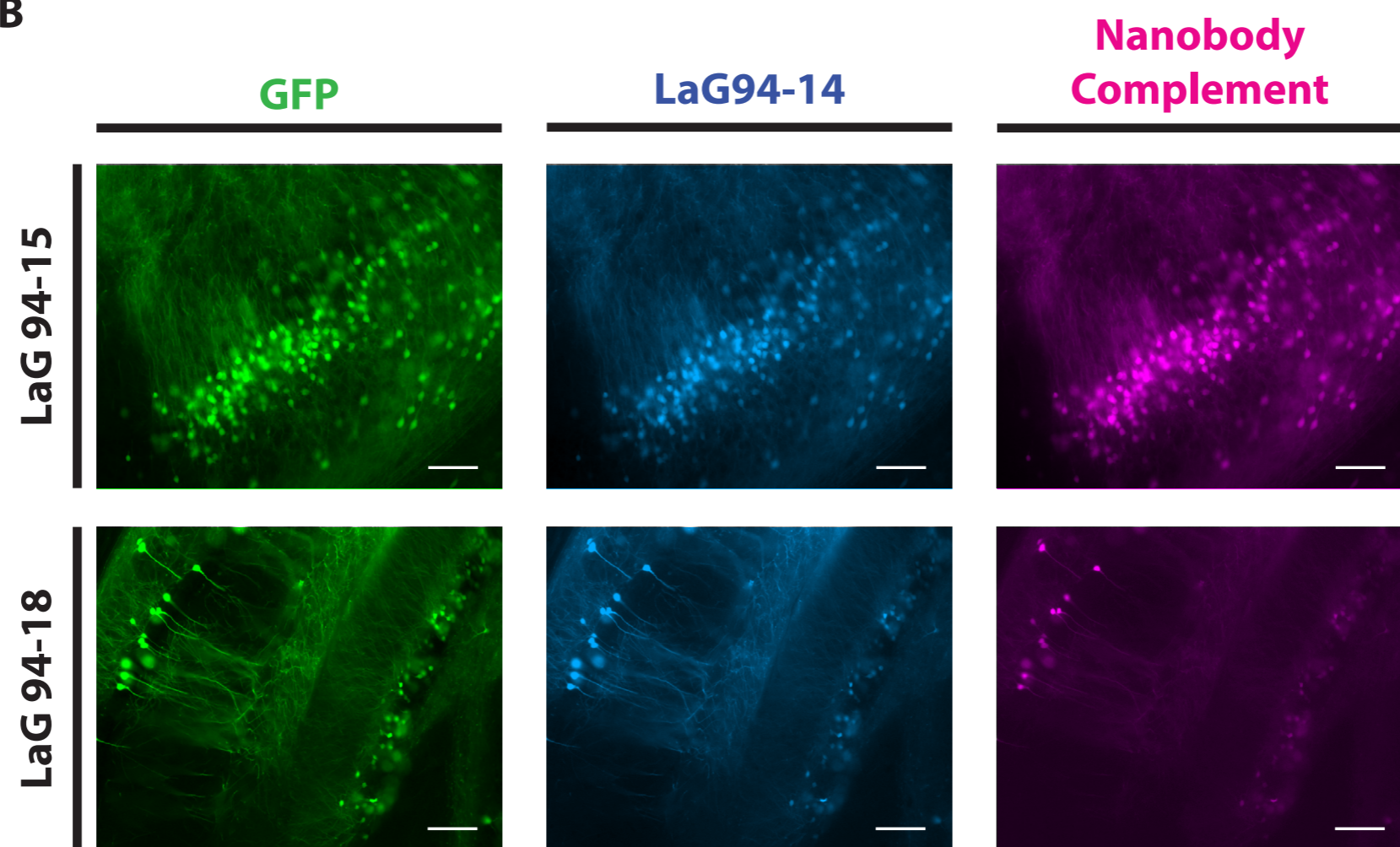


Figure S12

A



B



C

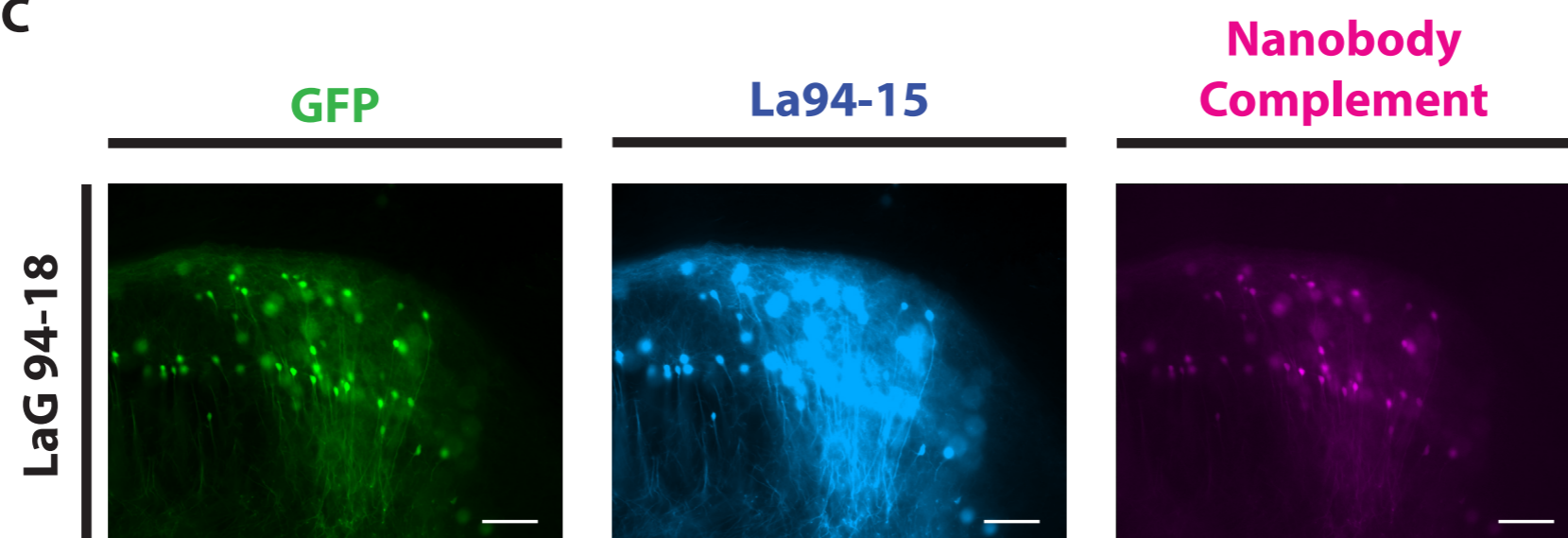
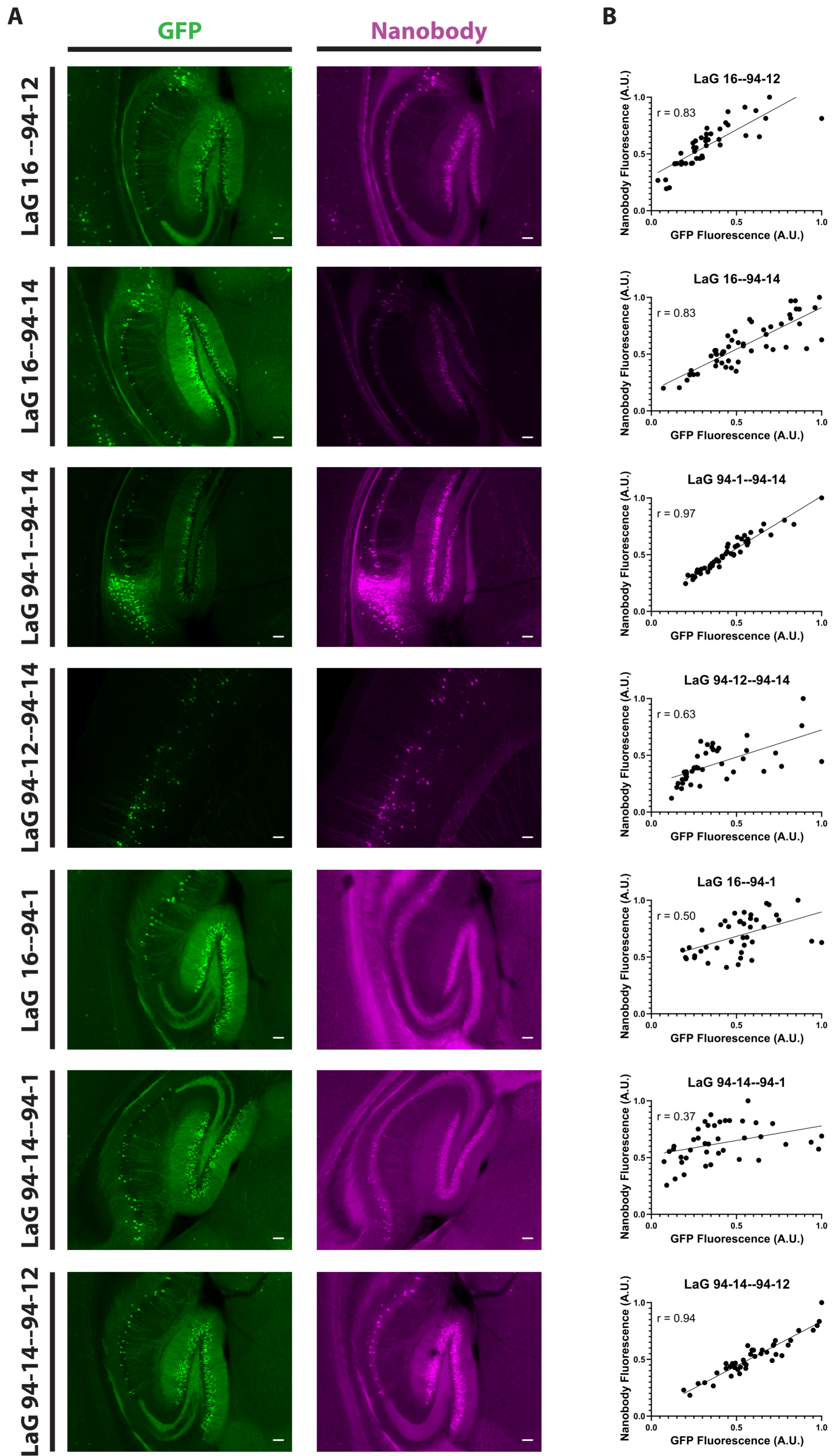


Figure S13



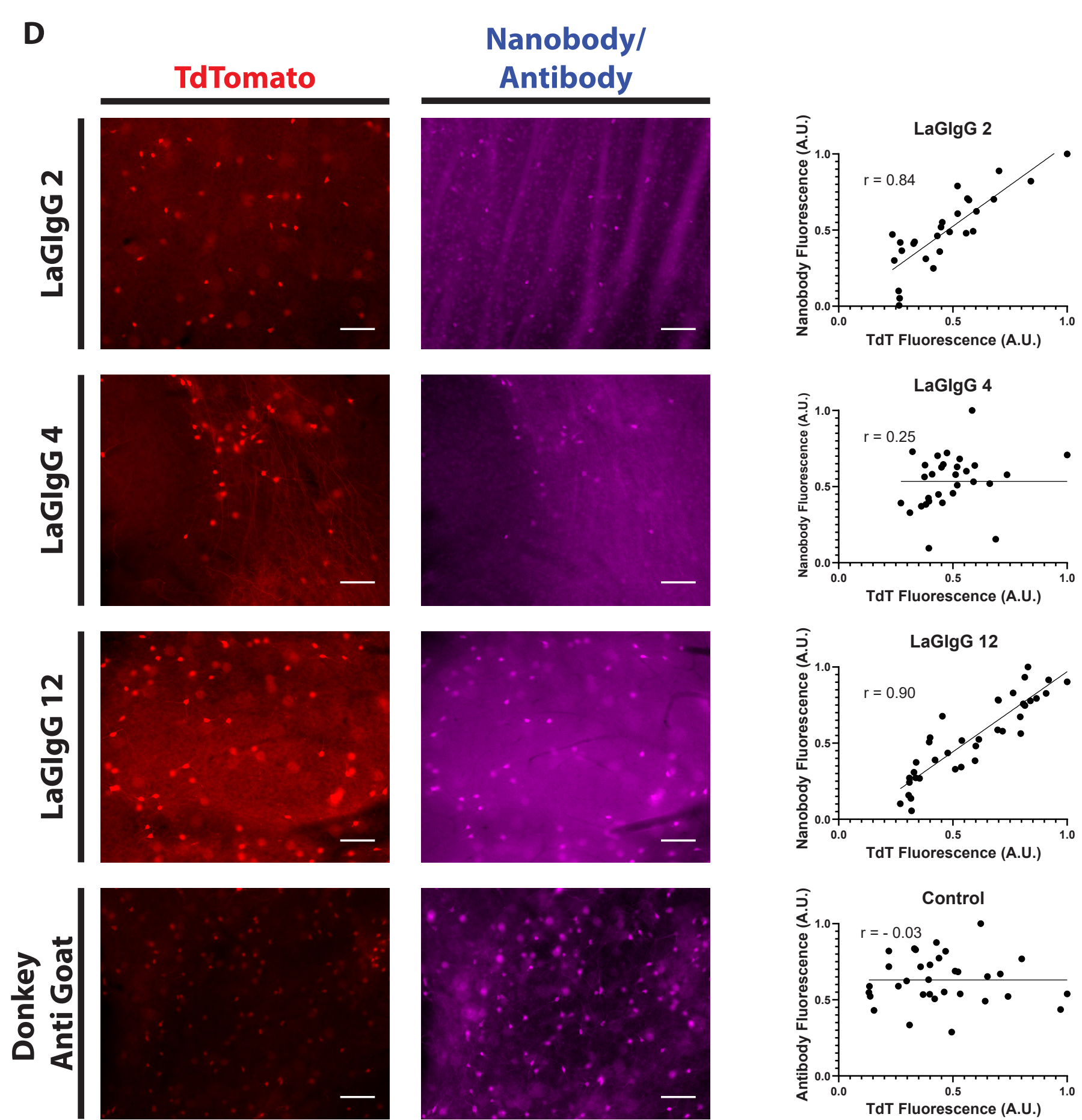
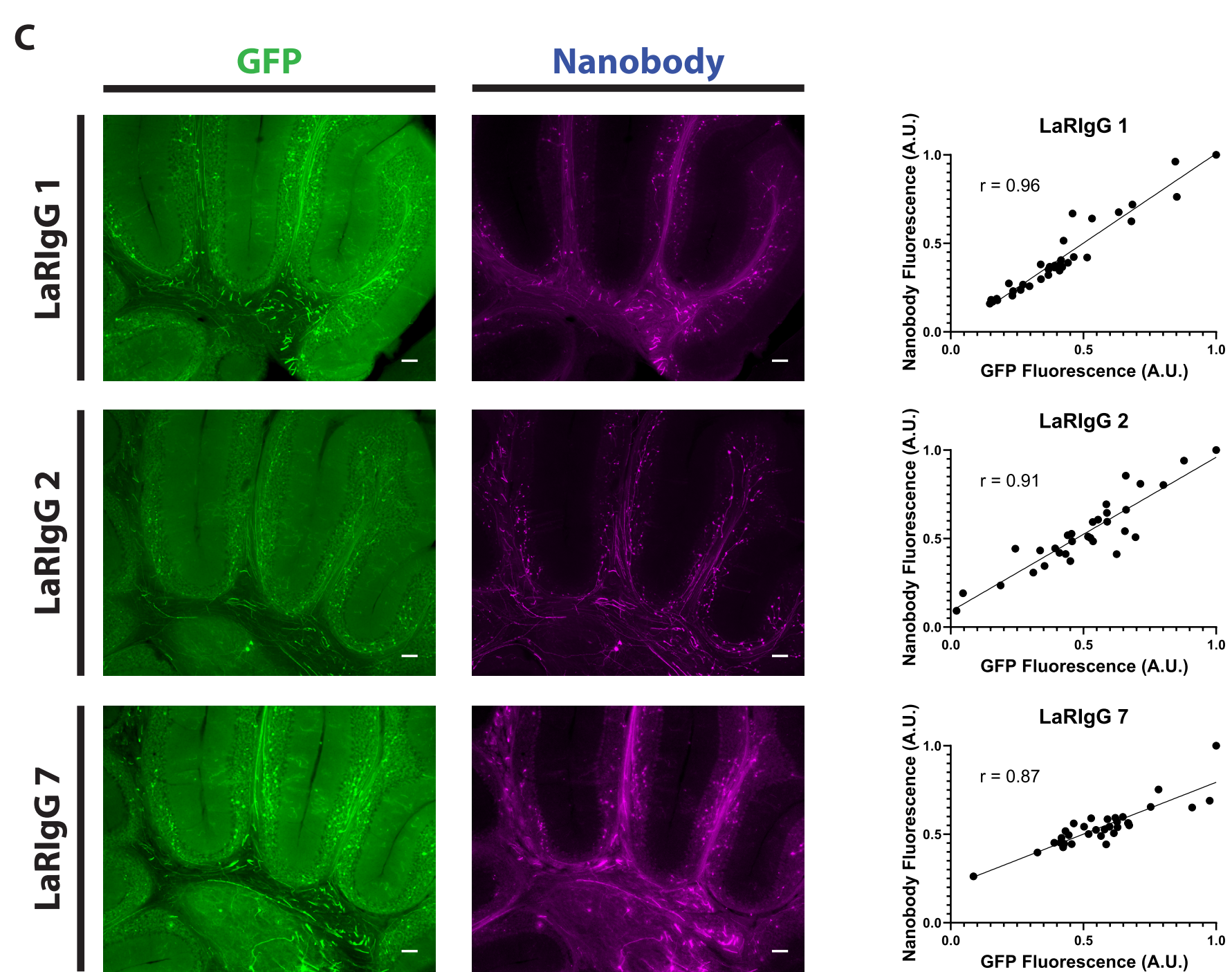
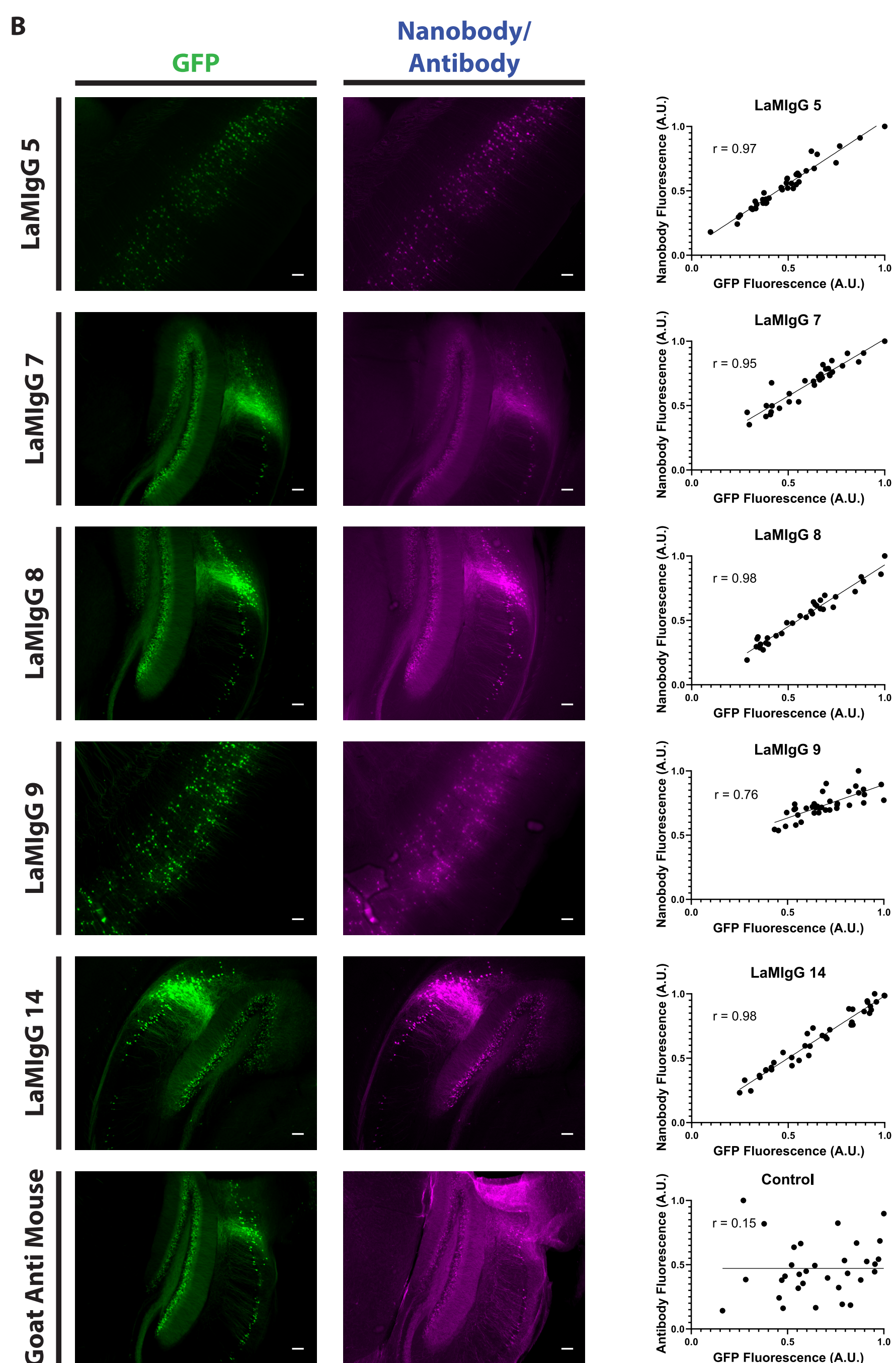
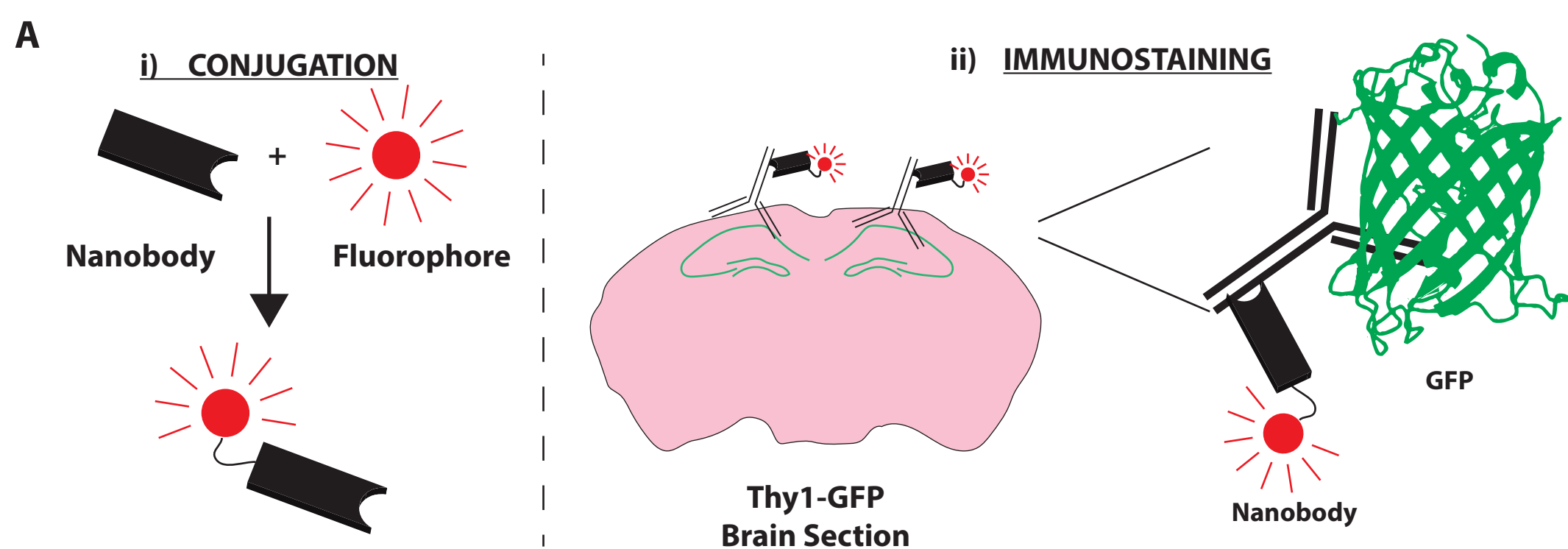


Figure S15

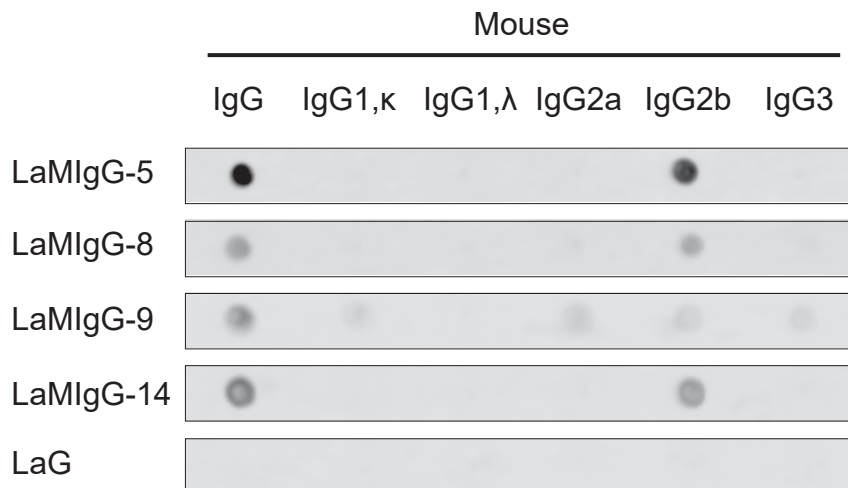
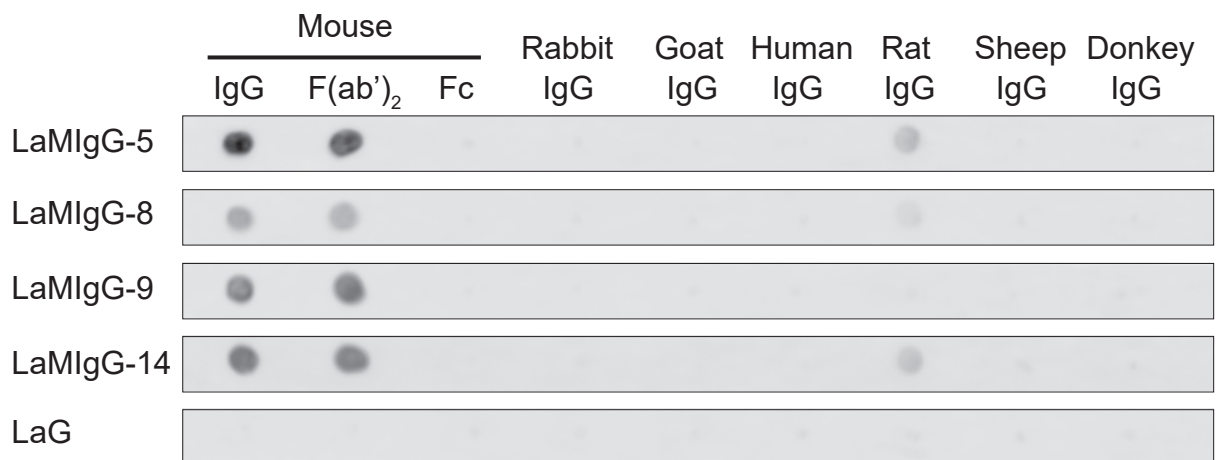
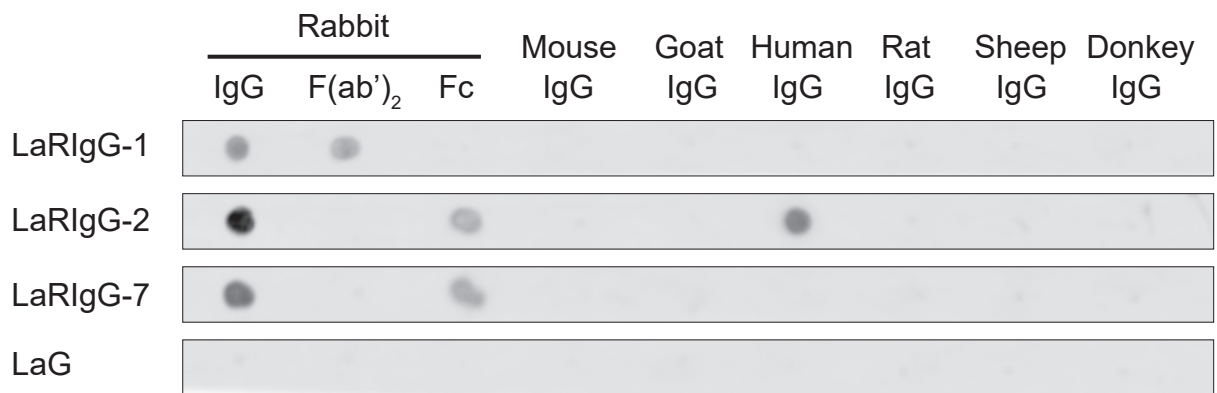
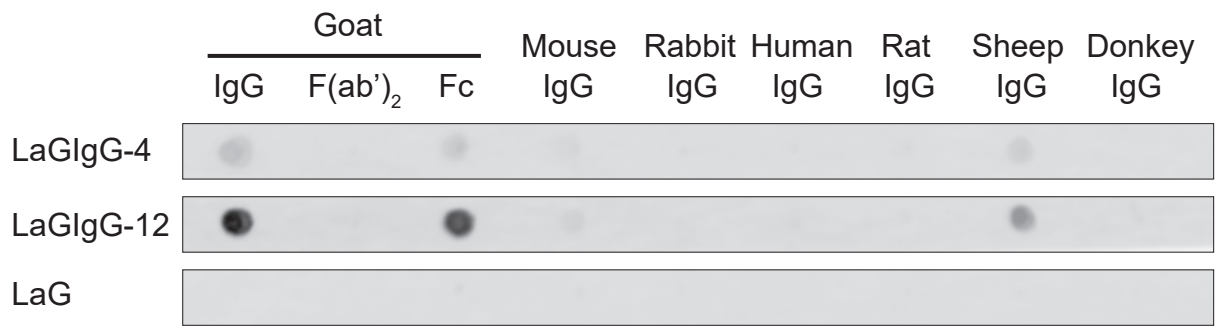


Table S1. Nanobody amino acid sequences.

Nanobody ID	Amino acid sequence
LaG94-1	MAQVQLVESGGGLVQTGDSLRLSCRASGLTLSVYSTGWFRQAPGKEHENV AWISRSGDETYEYEDSVKGRFTISRDNKNTVYLQMNDLKPEDMARYYCAAS NRGRTTAWNLYIYRGQGTQVTVS
LaG94-2	MADVQLVESGGALVQAGGSLRLSCEFSGGTFNLYGIGWFRQAPEQDHEFV AGVSRYGSTYESVSVKGRFTSSRDNAKNSAYLEMNTLKPEDTGVYTCAATV REQVPTNFGSYIYRGQGTQVTVS
LaG94-3	MAQVQLVESGGGLVQAGGSLRLSCTASGGTVSSSGMGWFRQSSGNDRDF VAFITPTGPTTWYEEVNSRFTVSRDNKNTIYLMKSLRPEDTAIYYCAAG NIGRDSRTYPYWGQGTQVTVS
LaG94-4	MADVQLVESGGGLVQAGGSLRLSCAASGDTSSIHFWEWYRQAPGKQRELV ATANNGDFTHLQDSVKGRFTISKDNAKSTVYLQMNDLKPEDTAVYLCYARSY SNSNFWGQGTQVTVS
LaG94-5	MAQVQLAESGGGLVQAGGSLRLSCAASGRTFAMGWFRRAPGKEREVVAS MSRSLDTHYADSVKGRFTISRDNALNTVYLQTDLSLKPEDTGDYFCAVGLR YWSLSGDYIRREKYDVWGQGTQVTVS
LaG94-6	MAQVQLVESGGGAVQAGGSLRLSCAASGLTFSLYTMGWFRQAPEKEPEFV AYISKSGDTHYADSVKGRFTIGRDNKSTVYLQMNSLKSSEDVAYYYCAASIK REITPNSIYRGQGTQVTVS
LaG94-8	MAQVQLVESGGGLVQAGGSLRLSCAASGRTFMNYAMGWFRQAPGKGRFV VASISWTGETSTYADFKGRFTISRDNKNTVYLQMDSLKPEDTAVYRCAAK RATDFDDTVREAAEADYWGQGTQVTVS
LaG94-9	MAQVQLVESGGGLVQPGGSLRLSCAASGGTFDHTMGWFRQAPGKQREF VSAIRLHGVGTYYADSVKGRFTISRDNKNTVHLMNMLKPEDTAVYYCAAD RNWSRATEKEEYDYWGQGTQVTVS
LaG94-10	MAQVQLVESGGGLVQAGGSLRLSCAASGLTLRVYDVTWFRQAPGKEREV GTVNRNGVWTDYADSVKGRFTISRDNKNTIYLMNSLKPEDTAVYYCAAKT RPKLSTLWDEYIYRGQGTQVTVS
LaG94-12	MAQVQLVESGGRLVQAGGSLRLSCAVSGGTFSLYSLGWFRQAPGKEREV AAFSREDGTTSFADSVKGRFTMFRDGTKNTVYLQMNSLQKEDTAVYSCAAT VRAQPSYSWSSYVYRGLGTQVTVS
LaG94-15	MAQVQLVESGGGAVQSGGSLSLSCAASGGIRISAMTWFRQAPGREREFV AAVTRTLGTTYTDSVKGRFTISRDSVKGLYLRMNMNMRPEDAAVYYCAAR TRGFSTKVALETTDGMWYWGKGLTVTVS
LaG94-18	MAQVQLVESGGGLVQAGGSLRLSCVASGRTFSTSAMAWWRQAPGKEREV VAAISWNVGTYYPDVSVKGRFTISRDKAQNAVYLMNSLKPEDTAVYYCAA RVTGYFGDQVRRADDYRYWGQGTQVTVS
LaTdT-1	MADVQLVESGGGLVPAGGSLRLSCAASTDLFSRYTMGWFRQAPGKEREV VSAITGMSSMTKYADSVKGRFTVSRDNKATMYLQMDSLKPEDTAVYYCAA HTTTAAATPQSMGLMLKYDHWGQGTQVTVS
LaTdT-2	MADVQLVESGGGLVQAGDSLRLSCAFSGSAFDTYAMGWFRQAPGKEREV VSAITWSSGNTYYADSAKGRFTISRDDAKNTVYLQMNSLKPEDAATYYCAA NRYGSSWSPLTKMGQYEHWGQGTQVTVS
LaTdT-8	MAQVQLVESGGGLVQAGDSLRLSCADAGLPYAFNHYGMAWFRQAPGKER EFVAGTSSGNTYYGDSVPDRFSFSRDNDKKTVHLHINSLKPEDTAVYYCAA RLGGRYEDASDYDVWGQGTQVTVS

LaTdT-9	MAQVQLVESGGGLVQP GD S L M L S C V I S G H T F T N F R M G W F R Q A P G K E R E F V A G I T W L F G R T D Y A D S V Q G R F T I S R D K G K D T M Y L Q M D S L K P E D T G V Y T C A V Q S G V L A Y P M R E G F Y D I W G Q G T H T V S
LaTdT-10	MADVQLVESGGRLVQDGGSLRLSCEVSGGTFDTFAMGWFRQAPGKEREF VSAITWSSGNTYYGDSVKGRFTISRNNVKRTVYLQMISLKPEDTGVYYCAA DPYGSSWSPLQKEGQYDYWGQGTQVTVS
LaTdT-38	MAQVQLVESGGGLVQAGGSLRLS CAASGRSFGNNMGWFRQVSGTEREF VAAISWVGLITEYSGAVKGRFTISRDN AKNTLYLQMNSLKPEDTAVYSCAAD GRDVQLPGRSEYDYNWYGQGTQVTVS
LaTdT-39	MAQVQLVESGGGLVQAGESLQLS CAASGGTLRDAIVGWFRQARGKEREF L ASITWVGLSTNYS DSVKGRFVISRANPQNTAYLQMNSLKPEDTAIYYCAAKI GAFGLADEFRTVPQYDSWGQGTQVTVS
LaTdT-44	MAQQVQLVESGGRLVQDGGSLRLSCEVSGGTFDTFAMGWFRQAPGKEREF FVSAITWSSGNTYYGDSVKGRFTISRNNVKRTVYLQMISLKPEDTGVYYCAA ADPYGSSWSPLQKEGQYDYWGQGTQVTVS
LaTdT-45	MTQVQLVESGGGLVQAGDSLRLS CAFSGSAFDYAMGWFRQAPGKEREF VSAITWSSGNTYYADSAKGRFTISRDDAKNTVYLQMNSLKPEDAATYYCAA NRYGSSWSPLTKMGQYEHWGQGTQVTVS
GST-1	QVQLVESGGGLVQAGGSLRLS CVCVGRAFGVYAMGWFRQSPGKEREFVA AITWNGGSIDYADSVKGRFAISRDNPPNTVYLQMNSLKPEDTAVYYCAKYRN NYYNEAPNSWGQGTQVTVS
GST-2	QVQLVESGGGLVQAGGSLRLS CVVVGRAFGAYAMGWFRQSPGKEREFVA AITWNGGSITYADSVKGRFAISRDNPPNTVYLEMNSLKPEDTAVYYCAKYTN NYYNEAPKSWGQGTQVTVS
GST-3	QVQLVESGGGLVQP GGS L R L S C A D S G R I F S D Y T M A W F R Q P P G K E R E F V A S I N R S G L G T Y Y A D S V K G R F T V T R D N A K N M V Y L Q M N S L T V E D T A V Y Y C G A R W G G S Y I D P R E G M Y D F R G Q G T Q V T V S
GST-4	QVQLVESGGGLVQAGGSLRLS CVCVGRAFGVYAMGWFRQSPGKEREFVA AITWNGGSIDYADSVKGRFAISRDNPPNTVYLQMNSLKPEDTAVYYCAKYRN NYYNEAPNSWGQGTQVTVS
GST-5	QVQLVESGGGLVQAGGSLRLS CVVTGRAFGAYAMGWFRQSPGKEREFVA AITWNGGTTNYADSVKGRFAISRDNPPNTVYLQMNSLKPEDTAVYYCAKYP NNFWHNEAPDSWGRGTQVTVS
GST-6	QVQLVESGGGLVQP GES L R L S C A A S G F T F S S H R M N W L R Q A P G K G L E W V S T T S T G G G S A A T Y Y V N S V K G R F T V S R D D A K N M L Y L Q M N S L K P E D T A R Y Y C A I P E S P G R P F G S T W S E Y R Y W G Q G T Q V T V S
GST-7	QVTLVESGGGLVQAGGSLGLSCTASGHTFNPYAMAWFRQAPGKEREFVAAI GWRVGTYYTDSVKGRFTISADNAKNTAYLQMNSLKPEDTAVYYCAAGFGR SPSNYDYWGQGTQVTVS
GST-8	QVQLVESGGDLVQAGGSLRLS CAASGSISAAYDMGWFRQVPGKQRELVAS L G S G G F S A Y A D S V K G R F T V S R D N A K N T V Y L Q M N S L K P E D T G V Y Y C A A R T P F D I W T T G R W G E Y D Y W G Q G T Q V T V S
LaMlgG-5	MADVQLVESGGGSVQAGGSLRLS CAASGFIFSKDVMNWFRQAPGKERELV AFINPSGTTTHYVDSVKGRFTISRDNVKSTVYLQMNDLKPDDTAVYYCQTRKF VGTTWRDYWGQGTQVTVS
LaMlgG-8	MAQVQLVESGGGLVQAGGSLKLSCATSGFTFSTSVMNWFRQAPGKERDM VAFINPSGSTNYVDSVKGRFTISR DSTKNTVYLQMNSLKPEDTAVYYCQTRL LVDSGWRDYWGQGTQVTVS

LaMlgG-9	MAQVQLVESGGGLVQAGGSLRLSCATSGRTGDLYAMAWFRQAPGAEREFV ATVTAIGGATYYADSVKGRFTISASYARNMVYLQMNSLKPEDSAVYYCAASM KRVWPLHRSDEAEYWGQGTQVTVS
LaMlgG-14	MAQVQLVESGGGLVQAGGSLRLSCAASGFTFSTSVMNWFRQAPGKEREM VAFINPSGSTHYVDSVKGRFTISRDSAANTVYLLQLNSLKPEDTAVYYCQTRL FVGAGWRDYWGQGTQVTVS
LaRlgG-1	MAQVQLAESGGGLVQPGGSLRLSCVSSGGTFGRYDMGWFRQAPGKEREF VAAISREVTNYGDSVRGRFTISRDNAAENTVHLLMNSLKPEDTAVYYCAADD AHRYSGTYYLSQDFRYTYWGQGTQVTVS
LaRlgG-2	MAQVQLVESGGGLVQAGGSLRLSCVASGGTFSDWPKGWFRQAPGKEREF VAAINWNGSGTTYADSVKGRFTISRDNAAENTVYLLMNSLKPEDTAVYSCAST FGIRSVGTNDYWGRGTQVTVS
LaRlgG-7	MADVQLVESGGGLVQTGGSLRLSCAASGRTFTTYDLAWFRQAPGKERECV AAVSRIGTTYYPDSVKGRFTISTDTAQNTVYLLQMNSLKPEDTAVYYCASAKF GSRWIRAKTSEGNYDDWGRGTQVTVS
LaGlgG-4	MADVQLVESGGGLVQAGGSLRLSCAASGITFSRSTMGWHRQAPGKRRELV ASINAVGSTNYVPSVKGRFTISRDNAAENTVTLQMNSLEPEDTGYYCVSDP YDDPRRYWGPQTQVTVS
LaGlgG-12	MAQVQLVESGGGLVQAGGSLRLSCSASGRTLSTYAMGWFRQAPGKEREF VAVISRSGGSTHYVESVKGRFTIVRDNAAENTVYLLQMDSLKPEDTAVYYCAVP KYTDYALSPEDFGSWGQGTQVTVS

Llama Magic 2.0

Detailed Methods

The Llama-Magic github repository is located at https://github.com/FenyoLab/Llama_Nanobodies. It includes all code, as well as detailed instructions for running the custom scripts for database preparation and for setting up and using the Llama-Magic web tool. Below is an outline describing the major steps in the pipeline.

Database Preparation

- 1) QC/Trim and Merge paired end sequencing reads: read pairs of length 300 bp are trimmed for quality and then merged. Trimmomatic (42) is used for trimming with parameters MAXINFO:275:0.8 (read 1) and MAXINFO:250:0.8 (read 2). Abyss (43) mergepairs is used for merging the reads with parameters -p 0.95 -m 25.

The following steps are performed using a series of custom perl scripts:

- 2) Convert merged FASTQ file to FASTA format.
- 3) Search for the PCR primers near the 5' and 3' ends of each merged sequence. This helps to determine the correct ORF for *in silico* translation, since the correct reading frame from the start position of the primer is known, and their orientation determines forward or reverse strandedness. This script adds the primer information to the header line of the FASTA file.

- 4) Perform an *in silico* translation to protein sequence. This script uses the output from step 3 to correctly perform the translation to protein sequence.

(The previous version of the script performs all possible ORF translations and selects the longest as the correctly translated protein sequence. This is a simpler method; however, it does not always result in the correct ORF.)

- 5) Perform an *in silico* digestion of protein sequences with trypsin/chymotrypsin. One or both digestions can be performed. Allowed missed cleavages = 1 (trypsin) and 4 (chymotrypsin). This script produces a FASTA file where each unique (chymo-)tryptic peptide is listed on a separate line. It also produces an index file for input to the peptide mapping step (step 8). This index file lists the sequence ID and all unique peptides associated with that sequence for each sequence in the database.

Llama-Magic: CDR finding, mapping of identified peptides and scoring of sequences

- 6) Run X!Tandem(44): the input to X!Tandem is the converted raw instrument files (mgf format) and the pre-digested database from step 5. The output is an XML file with identified peptides, expectation scores and other descriptive information. Detailed parameter settings for the X!Tandem run in the form of an XML file can be found at the github repository.
- 7) Parse the XML output file of X!Tandem to a tab-delimited text file with the following columns. Only matches with $\log(e) < -1$ are included.

<u>Column name</u>	<u>XML tag/attribute</u>
sequence	value of the seq attribute of the domain tag
log(e)	log10 of the value of the expect attribute of the domain tag
protein_uid	value of the uid attribute of the enclosing protein tag
domain_id	value of the id attribute of the domain tag
spectrum	text of the note tag with attribute label ="Description"

- 8) **Perform peptide mapping and ranking of proteins:** map the peptides from step 7 (peptides from the *in silico* digested protein sequence database that were identified by X!Tandem) to the protein sequences from step 4, and calculate a score for each sequence based on the location and number of identified peptides it contains. The score is defined as:

$$8 C_{CDR3} + 2 C_{CDR2} + 2 C_{CDR1} + 2 C_{SEQ} + \frac{1}{15} L_{CDR3} + \frac{1}{10} L_{CDR2} + \frac{1}{9} L_{CDR1}$$

$L = \text{length of region}$

$$C (\text{Coverage}) = \frac{\text{count of amino acids within an identified peptide}}{L}$$

The sequences are sorted by score and grouped by CDR3 region (proteins with CDR3 regions of hamming distance = 1 are grouped together). The output produced by this step is an interactive HTML file that allows the user to view detailed information on the peptide

mapping for each protein sequence. Each peptide that maps to a sequence is shown as a clickable link that will allow the user to view the detailed MS2 coverage for the identification. In addition, a uniqueness score is calculated as:

$$\text{Uniqueness} = \frac{\text{count of sequences in the group containing the peptide}}{\text{total number of sequences containing the peptide}} \times 100$$

If a user has performed both a trypsin and chymotrypsin digestion and has results from multiple MS/MS runs, it is possible to provide the results from both X!Tandem searches to this algorithm and it will combine the peptide identifications when mapping to the protein sequences and producing the sequence score.

CDR finding:

In order to score and rank the sequences, Llama-Magic implements a CDR finding algorithm as described in the table below. The listed positional indexing starts at the translated primer sequence starting position, which is where all the sequences start once they've been *in silico* translated (steps 3 & 4). For example, the indexing starts with the Valine at position 1 in this VHH sequence: 5'-**VLAAIGQG**VQAQVQWVESGGG...-3'. Llama-Magic clips the sequence at the 3' end following the sequence corresponding to the hinge region PCR primers: **PKTPKPQP** or **HHSEDP**. If the sequence of the short hinge is found, Llama-Magic will add 2 placeholder "amino acids"(**HHSEDPXX**) to keep the sequence lengths consistent for the CDR-finding algorithm.

The algorithm searches for certain conserved amino acid motifs at a range of positions to delineate the start and end of each CDR region as shown below. Alternate motifs are listed in subsequent lines and are searched if the previous motif is not found. When none

of the sequence motifs are found, default positions are used. Subsequent lines in the Default position column provide a last resort position for the case when neither the CDR start nor end position sequence motifs are found. The CDR3 region is first localized based on the ending position of the CDR2 region. If this fails, an earlier algorithm (from Llama-Magic 1.0) is used as an alternative (final 2 rows of table).

<u>CDR position</u>	<u>Search positions</u>	<u>Sequence motif</u>	<u>Position of CDR st/end</u>	<u>Default position</u>
CDR1 _{stpos}	30 to 36	SC C	C _{pos} +5	CDR1 _{endpos} -8 38
CDR1 _{endpos}	42 to 50	WXR W	W _{pos} -1	CDR1 _{stpos} +8 46
CDR2 _{stpos}	42 to 50	WXR W	W _{pos} +14	CDR2 _{endpos} -9 61
CDR2 _{endpos}	73 to 82	RF	R _{pos} -8	CDR2 _{stpos} +9 70
CDR3 _{stpos}	CDR2 _{endpos} to SEQ _{endpos}	DTAVYXC DXAVYXC DTXVYXC DTAXYXC DXXXYXC	C _{pos} +3	(CDR3 _{stpos} -alt)
CDR3 _{endpos}	CDR3 _{stpos} to SEQ _{endpos} CDR2 _{endpos} to SEQ _{endpos}	GTQVTV XTQVTV GXQVTV GTXVTV GTQXTV GTQVXV GTQVTX	G _{pos} -4	(CDR3 _{endpos} -alt)

CDR3 _{stpos-alt}	102 to 112	YXC C	C _{pos+3}	CDR3 _{endpos-11} SEQ _{endpos-27}
CDR3 _{endpos-alt}	SEQ _{endpos-20} to SEQ _{endpos-10}	WGQ WG W XGQ	W _{pos-1}	CDR3 _{stpos+11} SEQ _{endpos-16}